











level. The crucial technological ingredients in IRTS include low-cost sensors, clever software for analytics and visualization, and high computing power. Niculescu *et al.* [28] presented a smart infrastructure, which refers to all communication and information equipment that supports the development and the implementation of IRTS services. However, it mainly focused on smart rail infrastructure, with little consideration of smart rail communications. Above all, the discussions on smart rail systems in the existing literature have presented these systems as just some concepts the studies do not provide specific architectures or advanced technologies that could help in developing these systems.

### A. Network Architecture

We first review the general network architectures that have been proposed for 5G. Note that C-RAN has been designed and tested since 3G [29], and D-RAN have been applied to 4G as well [30]. Moreover, CUPS was already implemented in the 4G CN. However, due to the high demand for various services and due to the variety of communication scenarios, the performance requirements for wireless networks have become more stringent in recent years. To handle such requirements, 5G researchers have developed some novel RAN architectures, such as fog-computing-based radio access network (F-RAN), H-CRAN, and UDN. Based on SDN and NFV techniques, both RAN and CN can be made programmable, flexible, and elastic [15].

The HPNs (e.g., macro or micro BSs) existing in HetNets take charge of the seamless coverage and all of the control signaling interaction. Both cloud computing and deployment heterogeneity may ensure the backward compatibility to the current wireless system and for it to evolve smoothly to the next-generation networks. HPNs can be introduced into C-RAN systems to create H-CRAN, taking full advantages of both HetNets and C-RAN [31]. In contrast to the general C-RAN architecture, the H-CRAN system also supports the separation of the control plane (*C*-plane) and the user plane (*U*-plane). The paradigm HPNs deliver the overall control signaling throughout the system, including the signaling interaction with the BBU pool, such that the centralized control-cloud functions detached from the BBU pool can efficiently improve the round-trip latency. However, the constraints caused by the fronthaul links, as compared to those by ideal links, limit the overall SE and EE of the H-CRAN architecture.

To mitigate the shortcomings of C-RAN and H-CRAN, Peng *et al.* [32] proposed the F-RAN architecture. Incorporating fog computing into the edge cloud, which is a traditional RRU equipped with cooperative network functions (including web caching, signal processing, and radio resource management), can evolve into an innovative network element termed as fog-AP. Adjacent “smart” UEs (denoted as fog-UEs) can communicate with each other via D2D or relay modes to improve SE. The cross-tier interferences between the fog-AP and HPN can be suppressed in

**Table 1** Merits and Challenges of 5G Network Architectures

Network architecture	Merits	Challenges
H-CRAN	Backward compatibility, less handoff, front-haul efficiency	Serious interference in cross-tier and inter-tier
F-RAN	Front-haul efficiency, high capacity in hot spots	Complex and mass node reconstruction
UDN	Improve spectral and area throughput	Complex interference management
UCN	Reduce inter-user interference	Global CSI and huge fronthaul signaling

the BBU pool by using coordinated multipoint techniques. Four candidate transmission modes can be selected such that UEs can access F-RAN adaptively in accordance with the UE mobility characteristics: D2D and relay mode, local distributed coordination mode, global C-RAN mode, and HPN mode.

In contrast to previous cellular generations, 5G network infrastructure densification introduces the UDN paradigm to obtain a large system performance gain [33]. This is due to each user having one or more BSs exclusive service. Thus, proximal communications and improved spatial degree-of-freedom can be leveraged. However, UDN deployment also introduces highly challenging interference management. In order to tackle this challenge, efficient and realistic network deployment strategies are required by taking backhaul overhead, cost constraint, and computational complexity into account.

Another promising 5G network architecture is a UCN, which facilitates the centralized signal processing, low hardware cost, interference-free connection between users, and many more [34]. It takes advantages of the rapid development of cloud and edge computing techniques. However, UCN requires global CSI knowledge for efficient cooperation between APs. Furthermore, a large amount of fronthaul capacity is required for the UCN. In the future research, it is promising to investigate the dynamic AP selection, fronthaul compression, pilot assignment, power control, and mobility management.

Table 1 compares the merits and challenges of 5G network architectures.

With a huge available bandwidth in the mmWave band, such as the 60-GHz band and *E*-band (71–76 and 81–86 GHz), mmWave wireless communications can provide multigigabit transmission rates and support a lot

of high-speed data services [35]. Since small cells and macrocells are in different bands, there is no interference between small cells and macrocells. Combining the coverage and reliability from macrocells and high bandwidth from small cells, we expect small cells in high frequency bands underlying the macrocells to be able to deliver a good performance [35]. To ensure reliability and coverage, an intuitive method is to offload the control signaling to the overlaid MBS, providing continuous coverage and reliable connectivity at lower frequency bands, and the small cell provides high-speed data transmission at higher frequency bands within small coverage [35].

Compared with other electromagnetic waves at lower frequencies, mmWave communication has two main characteristics: high propagation loss and vulnerability to obstacles. MmWave directional BF is adopted to combat high propagation loss [36]. There have been several works on efficient BF for mmWave communications. By exploiting useful information such as location and mobility pattern, the overhead for BF is hugely reduced. There are several works on exploiting mmWave communications for HSR communication [37]. For example, prior works on mmWave channel modeling and mmWave BF design have validated the efficacy of HSR communications [38]. Furthermore, to combat the random blockage in practical scenarios, a variety of solutions have been proposed, such as relay-assistance [39], reflection exploitation [40], and large intelligent surfaces [41].

We now turn to the networks for railways. Considering the train safety-related services, mission-critical data information should be delivered with high reliability at lower frequency bands by the MBS, and thus only the passenger-oriented data information is offloaded to small cell BSs at high frequency bands [42]. The CUPS architecture is based on a HetNet deployment, in which small cells are overlaid on the coverage of macrocells. In general, macrocells working at low-frequency bands provide universal coverage and mobility support whereas small cells use high-frequency bands to provide high data rate transmission. Because frequency bands are scarce, passenger-oriented services requiring huge capacities depend on small cells for the  $U$ -plane. The corresponding  $C$ -plane bearing the control signaling remains in the macrocell and thus avoid frequent handover and severe intercell interference. By separating the user layer and the control layer resources, the operator can provide flexible services. The user layer can be distributed to different areas in order to get closer to the UEs, thereby reducing delay and bandwidth requirements. The control layer then reduces the complexity of operation and maintenance through centralized management [43].

Besides the frequent handover issue, another problem caused by handover is the group handover issue, as a number of active information users may desire to connect with Internet, enjoying multimedia services and online games. In this case, all active users have to launch their handover requests at almost the same time and all handover requests are required to be completed within tens of milliseconds,

which is challenging especially when the number of active users is large. This may put a heavy burden on both the train-ground communication and the handover computing processing. To alleviate such a problem, the two-hop architecture proposed in [44] may be an efficient solution. With the two-hop architecture, all information users in the train access the Internet via the train relay station (TRS) rather than directly linking the ground BS. In this way, the handover requests are aggregated and merged into a single one, so the UEs can be treated as a virtual big UE, and the handover processing burden is greatly reduced.

To guarantee the special needs of the required QoS, several studies have proposed to apply a network-slicing technology. This technology is a type of a virtual network architecture that allows multiple virtual networks operating on top of a common shared physical infrastructure to be separated in order to meet the specific needs of applications, services, devices, customers, or operators in terms of customized requirements such as latency, bandwidth, reliability, and security [45]. It is service-driven and aims at addressing the need of different use cases with highly diverse requirements. Network slicing further allows developers to set up new services and applications or modify existing ones. It can also support the communication services of a particular connection type with a specific way of handling  $C$ -plane and  $U$ -plane for these services.

Network slicing is created by E2E logical networks. Each network is flexible enough to provide one or more network services that are in accordance with the needs of the slice requirements such as vertical industry users, virtual operators, and business users [46]. Each slice can serve a specific vertical application to provide network services efficiently, to support the concurrent work of multiple use cases, and to provide flexibility to the network [47]. In other words, developers can take advantage of relevant technologies and key functions with network slicing; they can create specialized, dedicated, and logically independent virtual networks based on generic, programmable physical infrastructures. This technology also supports operators since it provides customized virtual networks for specific needs and scenarios and meets the KPIs. However, as far as we know, there is limited work on network slicing for railway communications.

An mmWave communication network architecture for railways has been presented in [48], emphasizing its capability to trigger handover in advance and thus enhance the probability of handover success. Yan and Fang [49] developed a network architecture that integrates mmWave communication and radar detection for railways. The proposed integrated network was deployed based on C-RAN. However, due to the difficult propagation characteristics of mmWave and due to the particularity of railway scenarios, there are still many challenges ahead. Moreover, these architectures are based on LTE but not on 5G.

A specific smart collaborative network architecture for railways was presented in [50], aiming to promote efficient and reliable communications under railway scenarios.

However, this architecture is not suitable for chain-shaped communication networks for the railways. It also lacks backward compatibility with redesigned network components and smart hierarchies.

## B. Channel Models

Channel models for railway communications are different from those for traditional cellular communications, which necessitates both new measurement campaigns and modeling approaches for the former. Wang *et al.* [51] provided a comprehensive review of the channel measurement campaigns conducted in different railway scenarios; another survey is provided in [52].

It is generally useful in channel modeling to distinguish between frequencies below 6 GHz and above 20 GHz. At below 6 GHz, most of the existing investigations have focused on narrowband measurements, characterizing path loss, small-scale fading, and shadowing in a variety of environments such as open area, cuttings, viaducts, and train stations. For example, Liu *et al.* [53] characterized narrowband railway scenarios at 2.35 GHz and then established a statistical position-based channel model. Based on the measurements performed along the Zheng-Xi railway line, a stochastic channel model for cuttings, viaducts, crossing bridges, and train stations at 930 MHz was provided by He *et al.* [54].

Many researchers have also used a variety of channel modeling approaches to characterize railway channels. Generally, geometrical approaches are preferred for characterizing the nonstationarities of the channels and to implicitly provide all the parameters required for multi-antenna characteristics. Ghazal *et al.* [55] developed a generic nonstationary wideband geometry-based stochastic model for MIMO systems in railway scenarios, and good agreement was obtained between the statistical properties of the proposed generic model and measurement data. Lin *et al.* [56] developed a finite-state Markov chain channel model while considering the impact of movement speeds on the temporal statistical channel characteristic. A geometry-based random-cluster model revealed the behavior of the clusters and their temporal changes for HSR [57]. Unlike in 3GPP-type models, these models characterize the time evolution of the clusters with high accuracy. Moreover, Zhou *et al.* [58] modeled the cluster evolution as well as the impact of railway-specific structures such as power poles based on extensive measurements. In summary, the channel characteristics and channel models under HSR scenarios below 6 GHz have been subject to a number of investigations but are far from completion and thus should be further explored.

In terms of mmWave bands, to the best of our knowledge, only measurements within a train car have been performed. Although no measurements between moving trains and infrastructure have been conducted, a few researchers have studied these channels through RT. Some of the environments that have been investigated include

a typical straight subway tunnel [59] as well as an arched tunnel [60]. The results have shown that systems with high gain directional antennas can support communication links longer than 1 km while inside a tunnel. The roughness of the wall inside the tunnel, antenna radiation pattern, and orientation in railway tunnels at mmWave bands have been studied in [61]. Meanwhile, Chang *et al.* [62] presented a channel model based on RT at 60 GHz for an HSR open area.

A 5G massive MIMO channel modeling method based on a theoretical nonstationary 3-D wideband twin-cluster channel model was presented in [63]. Moreover, the performance of 5G massive MIMO systems has been analyzed in [64] and [65], which reveals that the achievable SE can be assessed only through appropriate massive MIMO channel models. An important finding in [66]–[70] is that massive MIMO systems can be employed with low-cost hardware to achieve remarkable SE and EE. Guan *et al.* [71] defined and developed six scenario modules for mmWave and THz train-to-infrastructure channels for the “smart rail mobility.”

Overall, many articles have dealt with mmWave channel measurements and channel models; however, few have studied channel models at mmWave bands in rail traffic scenarios.

## C. Key Technologies

In railways, the requirement for seamless Gb/s-level transmission motivates researchers and developers to explore mmWave communication and BF technologies. However, an open and important question is whether mmWave and BF can really work in HSR. In order to answer this question, the researchers in Japan and South Korea have conducted field trials of mmWave communications for HSR to obtain high data rate [72]. Moreover, standardization efforts have been conducted on the mmWave-based HSR communications in 3GPP 5G NR specification from the aspects of network architecture, channel model and estimation, frame structure, Doppler compensation, efficient handover, and so on [73]. Therefore, mmWave is a promising and practical technique for HSR communications.

To extend the coverage in HSR, it is meaningful to use high directivity and narrow BF to combat the large isotropic path loss in the mmWave band. Several studies have dealt with massive MIMO for railways. Cui and Fang [74] designed a multistream BF scheme and used an adaptive beam-selection algorithm to exploit the train location information. However, the performance severely degrades when the interbeam interference increased.

Massive MIMO has been studied extensively as one of the core technologies in 5G to improve the system capacity [75]. SM is a candidate scheme of massive MIMO, effectively reducing energy consumption and complexity, making it a potential technique in beyond 5G. In particular, Cui and Fang [76] investigated the performance



of massive SM MIMO over a spatial-temporal correlated Rician fading channel for railway scenarios. The study theoretically determined that higher velocity makes temporal correlation more dominant when the train moves through a spatially variant field. Also, a hybrid analog-digital SM-BF scheme operating at mmWave frequencies for future railways was proposed in [77]. This scheme, however, performs poorly in the low SNR regime with velocities larger than 360 km/h. As such, Chen and Fan [78] suggested a low-complexity BF method based on the train location information, with a suboptimal solution for eliminating the interbeam interference and maximizing BS service ability. An mmWave BF scheme for railway disaster detection has been presented in [79]; in the study, the concerned area around the railways was divided into different detection areas with their corresponding danger sensitivity levels. In addition, the channel coherence time is small in mmWave communications with wide beams. Therefore, massive MIMO with BF is challenging in railway scenarios.

Channel estimation is of vital importance for wireless communication systems, especially in the case of high mobility. Fading channels for HSRs are fast time-varying due to train speeds of 360 km/h and more. Therefore, channel estimation for 5G on HSRs is a big challenge. Furthermore, channel estimation is a new research topic when combining with other 5G key technologies, such as massive MIMO and mmWave.

Currently, most of the existing articles that estimate the channel parameters for massive MIMO, both in TDD and FDD modes, assume that wireless channels are static or quasi-static. Besides, existing studies for channel estimation utilize training symbols and traditional estimators such as LS and LMMSE. However, as pointed out in [80], another two types of information can aid the channel estimation for HSRs: the wireless data of the past trains on the same track and the scenarios of the trains such as tunnels, viaducts, and cuttings. Accordingly, data-motivated and environment-sensing channel estimators, different from traditionally training-aided channel estimators, can be applied in 5G systems on HSRs to enhance the estimation accuracy.

As the Doppler spread increases in HSR, the resulting large intercarrier interference may degrade the performance of OFDM; thus, countermeasures or alternative modulation methods need to be developed. The orthogonal time-frequency space modulation, which performs modulation in the delay Doppler domain, has been proposed in [81] to address this issue. This method spreads each modulation symbols over the entire time-frequency domain to increase diversity.

As for the body of research that deals with reliable train-to-ground communication technologies, the most popular technique avoiding large penetration loss into the railway cars is the MR-based communication. UEs onboard the train communicate with the BS aided by antennas mounted on top of the train as gateways [82],

whereas distributed antennas constituting the other link end should be placed along the railway line. The key issue in the latter case lies on selecting an optimal cooperative remote antenna unit (RAU) of a linear distributed antenna for railway scenarios. Accordingly, Liu and Fan [44] presented a low-complexity approach based on the theory of convex optimization to select a transmit set of RAUs for capacity maximization. Likewise, Cheng *et al.* [83] proposed an adaptive antenna activation-based BF scheme to mitigate interbeam ambiguity. The proposed scheme can achieve SE gain over nonadaptive schemes and shows more robustness against direction-of-arrival estimation errors.

Moreover, Fan *et al.* [84] proposed a novel Doppler shift estimation algorithm for HST based on the reconstructed received signal via a discrete Fourier transform with massive linear receiving antennas. A filter bank multicarrier-based technique was presented in [85] to estimate the effect of time-varying channels on multicarrier signals, proving to be accurate for WiMAX and LTE at high speeds. High mobility also affects the media access control (MAC) and networking layer; thus, You *et al.* [86] presented a simple but effective distributed load-balancing algorithm to relieve service interruption due to frequent handovers in high-mobility scenarios.

As for the research on URLLC for railway scenarios, the first issue that needs to be considered is the design of the frame structure that would meet the requirements of ms-level latency. For example, the larger the block length is, the higher the coding efficiency would be [87]. However, this does not consider implementation complexity. For uncontrollable external interference signals in URLLC systems, Lee *et al.* [88] presented a “virtual pilot” frame structure that uses the data field to estimate the interference covariance matrix such that the transmission delay decreases. This frame structure, however, is suitable only for low-rate transmission scenarios; otherwise, the covariance matrix estimation error may cause an error diffusion problem, and consequently affect transmission reliability. As such, 3GPP has proposed the concept of mini-slots to meet the demands of URLLC [89] as it has been suggested that the subcarrier spacing would be larger under this condition. Variable length coding and feedback mechanisms can greatly improve the maximum achievable rate of the system [90].

Using time diversity via retransmissions can reduce the effect of deep fading. Accordingly, ARQ is the traditional retransmission technology, albeit reducing system transmission throughput. Although HARQ reduces the throughput loss, a major challenge in taking advantage of time diversity is how to characterize effectively the system interrupt probability after a certain number of transmissions [91]. If the code length is adaptively changed, then the system throughput can improve even if the HARQ feedback delay is larger. However, the delay introduced by the signaling overhead of ARQ mechanisms cannot be ignored. Therefore, some scholars have proposed

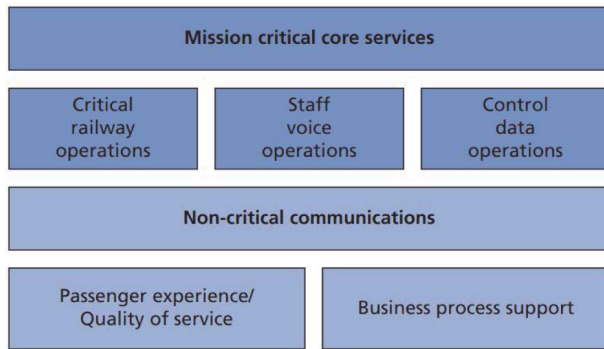


Fig. 2. Different categories of railway group services.

a method to avoid the retransmission of information confirmation.

The existing research on reliability lacks the design and optimization of the frame structure, retransmission, and feedback mechanism of URLLC for multiple users in railway scenarios. Thus, this subject needs to be further studied and explored. Most of the current frame structure designs are based on heuristics and can be evaluated only through numerical simulation. Therefore, researchers need to establish a novel optimization framework for designing URLLC frame structures while considering both signaling information and data payload. Furthermore, only a few studies exist with regard to the ARQ/HARQ mechanisms of short packet transmission in URLLC systems. The retransmission mechanism based on information acknowledgment makes the data transmission delay larger in URLLC systems, especially in cases of time-varying channel conditions.

User association was implemented using cell range extension by maximizing a user utility function by Hattab and Cabric [92]. Meanwhile, a scheme jointly considering user association and scheduling for load balancing in heterogeneous networks was investigated in [93]. Also, Su *et al.* [94] proposed distributed algorithms to solve the user association problem according to the backhaul capacity constraint and BF service in a two-tier heterogeneous networks. In addition, Kalantari *et al.* [95] presented an algorithm to find a suboptimal solution of user association to maximize the sum logarithmic rate of the users in a drone scenario. However, these schemes do not consider how to quickly switch associated BSs in a fast-moving railway network. Thus, it is a challenging task on how to deploy APs on the train or along the track to make a huge seamlessly connected heterogeneous network and establish a user association model for service provisioning in a heterogeneous wireless railway network with velocities larger than 360 km/h.

In railway environments, mission critical services, operation assistance services, and passenger services should be supported by reliable mobile communication systems, as depicted in Fig. 2. Mission critical services generally

include critical railway communications, train operational voice services, and operational data applications. Expected future railway mission critical systems should support intelligent transportation and control services, onboard and wayside HD video surveillance, distributed emergency communication, remote monitoring and diagnosis system, and so on. Due to the nature of the railway environment, future railway radio systems should fulfill the specific RAMS and QoS requirements demanded by railway services. In addition, mission critical services demand stronger delay, reliability, availability, and safety assurances, while additional services requirements are mainly based on the bandwidth capacity.

In general, there is limited work on network architecture, KPIs, channel models, and appropriate 5G key technologies for smart rail communication systems. However, these aspects are highly important to advance smart rail traffic systems.

#### IV. NETWORK ARCHITECTURE FOR 5G-R

The emerging services for smart rail may pose special requirements for the network architecture. Based on cloud and fog computing, we can expect that open network architectures will be converted from the seamless evolution of the current networks into SDN and NFV in order to facilitate CN slicing. The software-defined air interface, together with the anticipated advanced access technologies and physical layer technologies, can enable seamless radio access, optimize business management, and deliver diversified customer services and better QoE. However, the railway network architectures discussed in the existing literature focus only on TDMA and LTE-R that are not consistent with a 5G-R network architecture.

##### A. Proposed Network Architecture for Smart Railways

The requirements of the network architecture for railways are different from those of the public mobile communication network in many aspects. First, the network architecture for railways should be a dedicated communication transmission network that connects the dispatching center at all levels. This ensures that all elements interact with each other through different modes of information, such as voice calls between train driver and dispatching center, train operation control data, and video monitoring data for smart rail infrastructure. In addition, the rail communication network strives to establish an integrated and reliable emergency command system that would take real-time scene information as the decision-making base during emergencies (e.g., natural disasters or traffic accidents). Furthermore, the specific performance requirements of the train operation and control are embodied in reliable business, specialized equipment, timely transmission, and other aspects.

The biggest challenge for 5G-R is how to ensure the operational safety and reliability of the train and

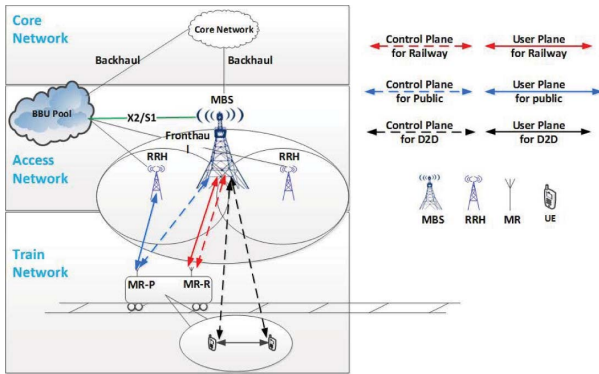


Fig. 3. Heterogeneous network architecture for smart railways.

the reliability of MBB communication for passengers at high-moving speeds and complex scenarios. In addition, the compatibility of legacy operational management systems should not be ignored. Accordingly, to protect the current infrastructure investments and to ensure normal operation during evolution, there should be no doubt about the long-term coexistence of multiple rail-oriented transportation systems. In other words, 5G-R will certainly evolve with the legacy networks, as it maintains backward compatibility and smooth transitions.

The design of a 5G-R network architecture mainly involves three parts: access network, CN, and air interface. Unlike C-RAN, D-RAN, F-RAN, and H-CRAN, the network architecture should fully take into account real-time and high-reliability access and always-online transmission. It should be capable of executing hybrid networking using both low (below 6 GHz) and high (mmWave) frequency bands. The CN design should likewise account for 5G-R services and applications, in which the network slicing and SDN/NFV networking architecture are used. Also, the 5G-R network architecture should support both ultrareliable train operation control signal transmission and high-quality user experience that would satisfy the diverse needs of diversified services. New CN elements should be added to provide different levels of QoS. For example, a train group call register needs to be added to the CN, and this should be different from the public mobile communication network.

There are two categories of MR: one that targets the passenger services (MR for passengers, MR-P), whereas the other focuses on train-related services (MR for railways, MR-R). In the heterogeneous network architecture for railways, as shown in Fig. 3, the MBS offers seamless handover and signaling interaction, the RRH provides data transmission, and the BBU pool performs signal processing and resource allocation. The MBS guarantees the reliability of critical missions, the compatibility with current BSs, and the wide-range seamless coverage.

Basically, there are three communication models for data transmission: one that the RRH delivers, the MBS, and another one that will be supported by

D2D communications. Unlike the situation of conventional BS deployment, a large number of the RRHs can be geographically distributed in a cost-effective and energy-efficient manner; thus, they are more suitable for remote deployment along the railway line. The BBU pool provides large-scale collaborative processing and enables tremendous amount of real-time connectivity. RRHs should be connected with the BBU through a high-capacity optical fiber fronthaul.

Fig. 3 shows three network domains that are based on the system resource layer, which connect the onboard access unit and the private network or the Internet. In the resource layer, the available physical resources can be virtualized into virtual resources. The radio access domain that is compatible with the current BSs integrates multiple access technologies, such as distributed heterogeneous hierarchical network, D2D connection, C-RAN connection, and satellite communications. We can expect to enable multitype and full-connectivity communications, allocate the system resource dynamically, and manage the network capability efficiently in railway scenarios.

The control domain based on the SDN and NFV establishes a componentized network. The network then decouples the *C*-plane and *U*-plane completely and flexibly reconstructs the network functions such as system control, network orchestration, and capacity openness. Accordingly, network slicing for specific use cases, reliability of train critical missions, and high QoE can all be achieved through the network function orchestration on the NFV. Meanwhile, the SDN enables the system to forward and control element separation completely. The forwarding domain concentrates on topology-aware routing and distributed forwarding of user data whereas the control domain facilitates the centralized control of the network. To meet the different demands of train-related and passenger-oriented services, the simplified design and flat distribution of the gateway equipment focus on traffic transmission and bandwidth expanding, enhancing both computation and storage capacity at the edge of network, and meeting the strict latency of time-sensitive service; therefore, they promote the full-capacity extension of the 5G-R network architecture.

In the logic view, CUPS has been clarified in railway scenarios. Subsequently, the specific network functions will be chained-up to deliver customized services, namely, the network slicing as a service. Fig. 4 depicts that the functional view consists of five layers and critical components. The modularized function layer is built on the infrastructure layer and plays a pivotal role. The service performance requirements lead to the demand for critical features, thereby providing network slicing as a service support for the upper applications. The MANO layer schedules and manages the system resources flexibly and collaboratively.

The infrastructure layer derives virtualized resources, accommodates the basic infrastructure, and facilitates the overall network service applications. Among these,

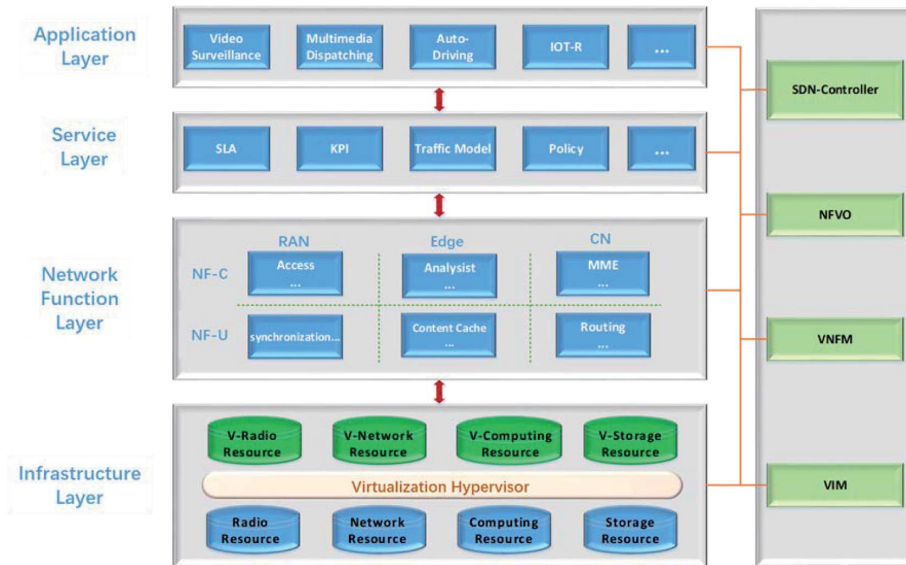


Fig. 4. Functional view of the 5G network architecture for railways.

the hypervisor virtualizes and manages physical resource allocation.

The function layer separates the *C*-plane and the *U*-plane, decouples the network functions from the hardware, constructs the modularized functions management system, and provides an E2E network-slicing framework. The network-slicing architecture encompasses the RAN, the transport network, and the CN resources or functions to exploit the multidimensional reorganization and association of function elements, thereby meeting the specific requirements of a certain use case.

The application layer provides diversified service applications to train access units and passenger terminals, consolidated application support and performance guarantee, and extra network capacity expansion and opening based on the customized individual network slices.

The service layer defines the features and requirements of diversified services, assuring the availability and reliability for all types of train-related subscribers. As such, this should be contracted between service providers and their customers. Specifically, the demanding description pattern comprises traffic characteristics (arrival rate, average packet size, and flow type); supplementary services (firewall service and open service chains); service-level agreement (customer technical support and traffic priority); and KPIs (RAMS, latency, and robustness).

The MANO layer consists of four function entities, namely, the SDN controller, the VIM, the NFV management entity, and the NFV orchestrator element. Among them, the SDN controller manages the network elements and controls the traffic processing. VIM deals with both physical and virtualized resources in the infrastructure layer. The VNF management entity supports VNF configuration and life cycle management [96]. Finally, the NFV orchestrator element arranges both the functions and resources.

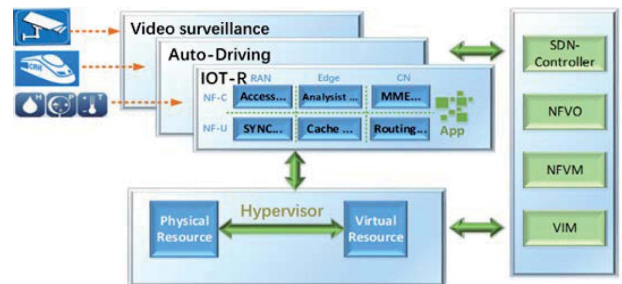


Fig. 5. Platform view of the 5G-R network architecture.

Four functional entities guarantee that the MANO layer fulfills the life cycle management and arrangement, thus implementing the E2E network slicing for railway wireless communications [97].

Fig. 5 shows the service and application platform for railway communications. E2E railway transportation network slices stretch across the business layer; the network and the infrastructure layer function through open interfaces by mainly adopting NFV/SDN for 5G-R to provide specialized business functions and customized protocol stack and orchestration. Context awareness allows the network to adapt to the needs of applications within the framework of network constraints and operator policy. In addition, the ability of both the network and device to use context awareness can help to further enhance user experience. This ability also enables the concept of the Internet to come to the user and provide the user with the most relevant and timely information, rather than the user having to go to the Internet to retrieve information and then filter out the irrelevant pieces of information. Context awareness includes awareness of the following.

- 1) Network analytics, including alternative radio access techniques, network layers (macrocell, mmWave, small cell, WiFi), and the corresponding congestion levels, capabilities, and performance characteristics.
- 2) Subscriber analytics, including subscription attributes, wireless activity level, loyalty management status, experience analytics, historical subscriber activity, location history, current location, subscriber contacts, and application usage.
- 3) Device attributes and capabilities, including information on single function versus multifunction devices, device support for specialized applications, machine-type communication versus subscriber devices, and radio and network optimization capabilities;
- 4) Application requirements, including QoS requirements, connection reliability, access price, power consumption, and security level.
- 5) Subscriber preferences, including preferred access options, power savings versus performance and access cost.
- 6) Operator policies and subscription context, including allowed services, service, attributes, and QoS.

Context information may be gathered from the device, network monitors, network elements, network databases, and analytics platforms. The network processes the context information when a device attaches to it or when an application is invoked. This determines the service attributes that govern how the network will treat the device and application. The service attributes for access may, for example, include cost, reliability, power consumption, security level, QoS, and mobility. The service attributes for access can be mapped to configurable 5G features, which the network assigns. For example, context information may determine that low-cost access, with no support for active mobility and long battery life, is best for providing service to a nomadic sensor device that attaches to the network. As a result, the network configures connectionless access with low priority, simple IP networking with no tunneling, and an idle-mode wake-up period of one day.

### B. Network Slicing for 5G-R

Network slicing is the key technology allowing to apply NFV to 5G. We can anticipate that a 5G-R network will be able to support a diverse set of emerging smart rail services and applications. Therefore, network slicing can support the multiple QoS needs of railways. Based on common infrastructures and on a range of techniques, network slicing can offer isolated and dedicated logical networks to certain users. However, the network slice types need to be further excavated and enriched.

In contrast, to our knowledge, network slicing for railway communications has yet to be investigated. In this section, we investigate how network slice types can be identified in terms of their service and use case characteristics. The proposed design process of the smart rail network slicing has the following procedures: network slice

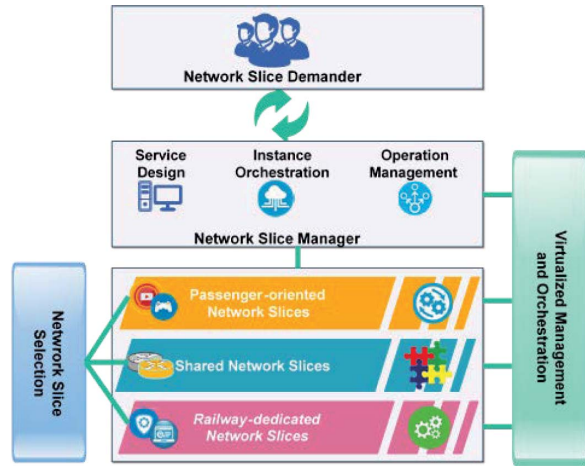


Fig. 6. Network slicing based on the architecture of 5G-R.

formation, network slice selection, network slice switching, user state maintenance, and new function identification.

As shown in Fig. 6, the network-slicing architecture for 5G-R contains two basic components: the slicing management and the slicing selection. The network-slicing management module is the core technology of the network-slicing operation maintenance. First, the user requirements at the service level are transmitted to the NFV/SDN-based cross-domain deployment descriptor and service profile, and then the whole life cycle of the network slice is managed to ensure the QoS requirements.

The network-slicing management function could provide an isolated, secure, and highly automatic E2E logical network to different kinds of service demanders, such as vertical industry users, virtual operators, and enterprise users. The process has three stages as shown in Fig. 6.

- 1) *Business design stage*: In accordance with the network slice templates and design tools, service demanders set the related parameters of the network slicing, including network topology, function component, interactive protocol, performance indicators, and hardware requirements.
- 2) *Instance orchestration stage*: A network slice management function sends the slice description file to the NFV MANO function module for the network slice instantiation. Furthermore, it will hand out the network element function configuration information via the interslice interface. Finally, it starts the connectivity test through migrating to active state.
  - a) *Implementation MANO*: An example of a management and choreography entity belonging to NFV; NFV MANO is broken up into three functional blocks.
    - i) *NFV orchestrator*: Responsible for onboarding new network services and VNF packages, network service life cycle management, global resource management, and validation and

authorization of NFV infrastructure resource requests.

- ii) *VNF manager*: Oversees the life cycle management of VNF instances and the coordination and adaptation role for configuration and event reporting between NFV and E/NMS.
  - iii) *VIM*: Controls and manages the NFV infrastructure compute, storage, and network resources.
- 3) *Operation management stage*: While in a running state, the slice owner can perform real-time monitoring and dynamic maintenance of the slice with relevant slicing management functions including dynamic allocation of resources, design, creation, deletion, and update of the network slices as well as the fault-alarming module.

The network-slicing selection module maps the user terminal to the appropriate network slice. Considering services subscription and functional properties and other various factors, it is capable of provisioning proper network-slicing for different user requirements. Shared or separate network slices can be accessed by different users. Also, the slicing selection contains RAN slicing (LTE or wired networks) and CN slicing (various combinations of *C*-plane and *U*-planes from different CNs). Industrial network slice instances include intelligent control system for railways (in need of low latency and high reliability). This control system is required for dedicated RAN and specific CN and for broadband communications, which is required for proper public mobile network in cooperation with trackside private network for railways.

E2E network slices connect the UE, access slices, CN slices, and the mapping relationship between the slices into a complete logical network. The slicing selection function routes data transmission from RAN slices to appropriate CN slices. The CN slices compose a set of supporting service functions that can either be shared or be exclusive to certain E2E network slices. The virtual LAN technology of virtual-switch can enable service isolation inside a network slice whereas the virtual-switch technology of the virtual-router can offer the business isolation between network slices.

*Service slicing category*: The nine types of network slicing are summarized as follows.

- 1) *Auto-driving slicing*: Shows strict requirements for high safety, high reliability, and low latency.
- 2) *D2D slicing*: Establishes direct communication links between the interactive devices and enables low-latency and large-capacity data transmission.
- 3) *MBB slicing*: Provides travelers with access to MBB communication networks and supports interactions inside the coaches.
- 4) *Ultrahigh reliable slicing*: Attends to those services that need ultrahigh reliability.
- 5) *Ultrahigh available slicing*: Fits with applications that call for ultrahigh availability.

- 6) *Ultra-HD slicing*: Capable of ultra-HD video streams.

Meanwhile, the three types of IoT slicing are as follows.

- 1) *Massive IoT*: Monitors the status of the massive number of static sensors of the fundamental fixed infrastructure and of the train equipment and requires high reliability.
- 2) *Scheduling IoT*: Delivers the information of train scheduling and crew arrangement with high flexibility.
- 3) *Asset IoT*: Responsible for passengers luggage, cargos, and containers and also supports railway ticketing.

## V. CHANNEL MODELS AND OPERATING FREQUENCY BANDS FOR SMART RAIL

The characteristics and models of wireless channels for smart rail are critical for evaluating 5G-R key technologies. Among other decisions, the operating frequency bands suitable for rail communication systems are determined based on the channel models. In this section, we investigate the channel models and operating frequency bands for smart rail.

### A. 5G Physical Communication Scenarios for Smart Rail

Due to the variety of applications and deployment scenarios, 5G networks need to support extremely varied requirements and propagation channels with different properties and attributes. For example, the path loss and multipath effects are different in different propagation environments. Thus, radio wave propagation scenario partitioning is a fundamental step for wireless channel modeling. Although the International Mobile Telecommunication System-2000 [98], the Universal Mobile Telecommunications System developed by 3GPP [99], and the WINNER project group [100] all define typical scenarios, none of them specifically relate to railways. On the other hand, railway scenarios (e.g., cuttings, viaducts, tunnels, and crossing bridges) significantly affect the propagation characteristics. Therefore, developers need to establish a detailed scene-partitioning scheme for railways to improve the quality of wireless network planning and optimization. Based on practical measurements of typical railway lines and railway stations, and considering some other factors, such as physical and user attributes, we have defined 16 scenarios for railways [101].

Another important aspect in evaluating 5G-R key technologies are the communications scenarios, which mainly include train-to-infrastructure, intercarriage, intracarriage, inside station, infrastructure-to-infrastructure, and train-to-train.

- 1) *Train-to-infrastructure*: The links are between the AP/transceivers of the train and the infrastructure nodes of the fixed networks. The links deliver bidirectional streams with high data rates and low latencies and provide robust communication links with

latencies lower than 10 ms, together with an availability of 99.9%–99.999%, while moving at speeds of up to 500 km/h.

- 2) *Intercarriage*: A wireless network operates in-between carriages since it is very costly to wire a train for network access. Likewise, rewiring when a train needs to be reconfigured is inconvenient. This scenario of a wireless network requires a high data rate and low latency because the APs are arranged in each carriage such that each AP serves as a client station for the AP in the previous car while also serving as an AP for all UEs within its car.
- 3) *Intracarriage*: The links provide wireless access between the APs in the carriage and the passengers or the sensors of the equipment inside the carriage. In this scenario, real-time HD videos need to be accessed with low latencies.
- 4) *Inside station*: The links provide wireless access between the APs and the UEs in railway stations. Passengers can have access to MBB communication services. Furthermore, stations provide communication infrastructure to support commercial (e.g., cash desks) and operational services (e.g., automatic doors, surveillance, and fire protection).
- 5) *Infrastructure-to-infrastructure*: HD video and other information is transmitted in real time between multiple HDTV IP/HD-SDI cameras. The APs are deployed at station platforms and at the wayside along rail tracks as a high data rate WB or as IoT. The infrastructure nodes are connected in real time and interactive; they are supported by bidirectional data streams with very high data rate and low latencies.

## B. Channel Characteristics and Models Below 6 GHz in Smart-Rail Scenarios

In recent years, researchers and developers have conducted extensive sub-6-GHz channel measurements in various railway environments. Several studies on railway propagation channels focus on narrowband fading and the behavior of composite channel parameters. For example, the analytical model in [61] provides a fit for path loss, shadowing,  $K$ -factor, and other statistical channel parameters. Meanwhile, Guan *et al.* [102] and He *et al.* [103] studied the impact of viaducts, cuttings, stations, and cross-bridges on path loss and fading. Based on the extensive channel measurements conducted along the Zhengxi railway line, the PLEs are estimated, for example, PLE of the viaduct scenarios is about 3.5, and the PLE of the cutting scenarios is around 4.3. Meanwhile, He *et al.* [104] proposed a path loss model at 930 MHz based on Hata's formula that is modified with two correction factors to ensure a sufficient fit in different environments.

For the amplitude distribution of small-scale fading in railway scenarios, in most cases, the Rician distribution provides the best fit [67]. Therefore, the Rician  $K$ -factor can be used as a useful parameter in assessing the severity of fading and the resulting impact on system reliability.

Measurements show that one can model the dB-scale  $K$ -factor as a Gaussian random variable that has mean  $m_k$  and STD  $\sigma_s$  that are dependent on the location within the cell. The parameters of the distributions are summarized in [103].

## C. Channel Characteristics and Models at mmWave Bands in Smart Rail Scenarios

Recently, researchers have gradually shifted their focus from railway channel measurements to ultrahigh frequency (mmWave) range, emphasizing multidimension channel characteristics, for example, delay, Doppler, and spatial domain.

Prototype experiments on vehicular environments have confirmed the viability of highly mobile mmWave communications. For example, Samsung achieved in 2014 a rate of 1.2 Gb/s on a vehicle-to-BS communication link at 110 km/h in the 28-GHz band [105]. In 2015, Korea Electronics and the Telecommunications Research Institute achieved 1-Gb/s transmission rate in the 30-GHz band in the Seoul subway [106]. However, due to the lack of channel measurement data, simulation technology, and channel models, neither academia nor industry can fully evaluate and optimize their system performance under various HSR communication scenarios at higher speeds. The challenges include the exploration of static or dynamic radio propagation mechanism and multipath birth and death mechanisms, mmWave large-scale antenna channel characteristics, and a modeling theory for HSR scenarios.

Yang *et al.* [107] and He *et al.* [108] measured the mmWave propagation in railway scenarios, and accordingly studied the path loss characteristics. However, the measurement data do not provide a full picture. Moreover, the measurement data validate only the path loss model while the other channel characteristics are not verified since relevant data were lacking.

Other effects that have been mostly ignored up to now and thus will require particular attention are as follows: 1) the elevation characteristics of the propagation (since radio waves propagate in three dimensions and reflection/scattering disperses radiation in azimuth and elevation) and polarization and 2) 3-D channel measurements are necessary, especially in designing and evaluating smart antenna and massive MIMO. Similarly, polarization needs to be taken into account to assess its impact on the possible diversity and multiplexing gain.

Existing mmWave directional channel measurements were not performed under dynamic environments due to a lack of measurement equipment. A few exceptions exist, but none of them were done in railway environments. There are no measurement results under such circumstances since it is difficult to build suitable channel sounders for railway mmWave applications. Hence, in railway scenarios, comprehensive mmWave channel measurements are needed to establish and/or validate suitable channel models.

**Table 2** Definition of Six Modules for 5G mmWave Railway Channels

Module index	Module 1	Module 2	Module 3	Module 4	Module 5	Module 6
Definition	Tunnel entrance on steep wall connecting the cutting to the crossing bridges	Viaduct with open train station	Urban with semi-closed train station	Rural with CCT tunnel	Rural connecting double-track tunnel	Single-track viaduct
Special objects	Steep wall, cutting walls, crossing bridges	Open train station, indicators	Buildings, semi-closed station	Vegetation, CCTs	Dual-track tunnel, barriers	Single-track viaduct, barriers
Common objects	Trains (metal), tracks (concrete), pylons (metal), traffic signs (metal), billboards (metal/LED/concrete)					

To verify new communication regimes in railway environments, we first need to define 5G mmWave railway scenario modules together with their distinct propagation features. Six modules, representing various typical railway scenarios, were constructed in [109]. Note that the coverage ranges for which these modules are valid range from 100 to 200 m, which is much shorter than the macrocell or microcell range with carrier frequencies below 6 GHz. The high isotropic path loss of mmWave bands, which is partially compensated by BF gain, considerably shortens the link length, thereby providing us with the chance to concentrate on the six modules in order to compose a comprehensive smart rail environment. As summarized in Table 2, apart from common objects (e.g., trains, tracks, pylons, traffic signs, and billboards), each module includes its special objects, that is, steep walls, cutting walls, crossing bridges, train stations, indicators, barriers, vegetation, CCT, dual-track tunnel, and single-track viaduct. The 3-D models of both the comprehensive scenario and the six modules are publicly available and freely downloadable (see <http://raytracer.cloud>).

**Module 1:** Tunnel entrance on steep wall connecting cutting with crossing bridges: In a railway environment, a tunnel is an artificial underground passage that is specifically built through a mountain. Thus, when the mountain is high, rail builders usually construct steep walls to protect the tunnel entrance. Such steep walls are huge reflectors for wave propagation. Cutting is another common scenario in railways. A rail cutting is a man-made valley that carries the track as its base; it is used to pass through the hill that is not particularly high. Cuttings are made of concrete and stone with some vegetation on the surface. The depth of the cutting is usually 3–10 m. The slope height is about 13–14 m, and the inclination angle is 35°–40°. The distance between the two bottoms of the slopes is 12–13 m. The most common semi-closed obstacles in cutting walls are crossing bridges. They can block the LoS connection and generate strong multipath propagation.

**Module 2:** Viaduct with open train station: A viaduct is a long bridge across the uneven ground found in rural or urban settings. It is a common scenario in railways; for example, viaducts make up almost 70% of the Beijing–Shanghai railway line. We assume that we can obtain a LoS propagation condition most of the time. This module is composed of an open train station (awnings that only cover the platforms and not the rails) and a viaduct in a rural environment. We note that there are three types of modern train stations: closed type, semi-closed type, and open type. The first type is similar to the typical indoor environment whereas the latter two types are more particular. Usually, smaller, open stations appear in rural scenarios whereas larger, semi-closed stations are constructed in urban environments. Thus, this module is composed of an open train station and a viaduct in rural areas to represent the features of these two structures.

**Module 3:** Urban with semi-closed train station: The urban scenario represents the propagation in urban areas when the train is entering, leaving, or passing the city. Most of the buildings are higher than 10 m and are at least 10 m away from the barriers. In the semi-closed scenario, the huge awnings are usually designed to keep the rain from the passengers and trains. Thus, the channel will exhibit strong multipath because the semi-closed station, barriers, and buildings are present in this module.

**Module 4:** Rural with CCT tunnel: The rural scenario represents the environment where there is a large range of open areas, very few buildings, and where certain vegetation grows adjacent to the track. A CCT tunnel is a method for building shallow or short tunnels, in which the builders excavate a trench and then roof it over with an overhead support system. The overhead support system should be strong enough to carry the load of what needs to be built above the tunnel. CCTs are built in trains for two reasons: 1) to prevent potential landslides and 2) to cover the rail as a defense against wind and other unfavorable climate conditions. In fact, a CCT is not a complex construction method; however, it is challenging for wireless



planners and system designers of railway communications. Thus, in this module, the CCT, barriers, and the dense vegetation adjacent to the barriers are the main structural and influential factors that need to be evaluated.

**Module 5: Rural environment connecting double-track tunnels:** The cross section of double-track railway tunnels is usually vaulted or built in a semicircle fashion, with a height of 5–10 m and a width of 10–20 m. Two trains have the chance to pass by each other inside the tunnels, and the effect of this situation needs to be further evaluated. Thus, we design this module such that we can study the channels that connect a double-track tunnel in rural settings.

**Module 6: Single-track viaduct:** The difference between two single-track viaducts and one double-track viaduct is the presence of barriers between the two tracks. In this module, we can further determine the influences of barriers and railway crossing. However, the barriers and tall objects (e.g., tall trees or billboards) on both sides of the viaduct still affect propagation. The train station is another kind of important structure along the railways. The huge awnings, LED indicators, steel frames beside and above the track, and metallic pylons can block parts of the LoS, and consequently influence the channel. Note that if the target cell includes only a part of the scenarios, the users can further divide this module into more submodules, such as pure viaduct and pure open train station.

MmWaves are highly sensitive to the propagation environment and mobility. As such, multiple antenna technologies adopted in the mmWave band, such as BF, have very strict requirements for the spatial resolution of the channel. As discussed above, it is difficult, expensive, and complex to perform dynamic channel measurement at mmWave bands. This makes it impossible for us to rely only on measurements in order to obtain comprehensive and elaborate channel properties (e.g., time, delay, and frequency domains). Accordingly, the design of BF algorithms and other system design questions can be investigated based on detailed simulations of electromagnetic propagation.

The most suitable way to simulate radio wave propagation at high-frequency bands is through RT, which approximates electromagnetic waves as quasioptical rays. Propagation at high frequencies is similar to optical rays; thus, an RT approach is inherently suited to obtain highly accurate channel simulation and modeling at mmWave or even THz frequency bands. Different rays that correspond to different multipath components (MPCs) are traced from the Tx to the Rx. Each MPC is characterized by multiple output parameters: time delay; real part of field intensity; imaginary part of field intensity; AoD azimuth; AoD elevation; AoA azimuth; AoA elevation; type (transmission, reflection, scattering, and diffraction); and number of reflections or other interactions.

Although RT is well suited for high frequencies, both the calculation complexity and the simulation time exponentially increase as the number of objects and faces in a 3-D-scenario increase. Yet the required resolution

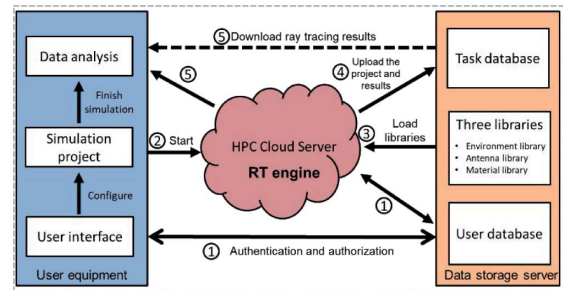


Fig. 7. RT simulator based on HPC.

(number of faces representing complex objects) increases with frequency. This makes it much more challenging to achieve accurate and efficient RT simulations for railway communications at mmWave and THz bands than when simulating at traditional frequencies below 6 GHz.

One promising solution to address this challenge is to transfer an RT simulator from a personal computer to HPC platforms. This can make RT much more efficient without sacrificing accuracy. Thus, we develop a high-performance RT supercomputing platform. All the components are connected via the Internet such that the components can exchange commands and data. As shown in Fig. 7, the platform is composed of five layers: data input layer, data transmission layer, data storage layer, data analysis layer, and application layer. Developing the smart rail RT includes building the RT database and material parameters, RT correction, antenna modeling, RT calibration, and channel modeling based on the RT simulation.

To build the RT database, we first need to parametrize the target smart rail scenarios. We then get the typical objects and materials from the scenario model. According to the parameters (electromagnetic parameters, roughness, etc.) of the individual models of radio propagation mechanisms, (i.e., LoS, reflection, scattering, transmission, and diffraction), we need to construct the material library for smart rail systems (frequency band up to 100 GHz). For different antenna types (single antenna, multiple antennas, and large-scale multi-antenna) and different array types (linear array, planar array, cylindrical array, etc.) that work in different frequency bands and operation mode, we need to study the antenna modeling approach and then establish the antenna library. The material library can be calibrated by comparing the RT results and sample measurements; the latter do not need to be performed at full train speed but can be performed quasistatically. We then use the SAGE algorithm to extract MPC parameters from the measurements and then calibrate the RT toolbox based on the reflection mechanism [109].

The data input layer mainly collects various data including environment information, antenna patterns, Tx and Rx deployments, frequency, and channel measurements, and accordingly transforms them into predefined structures. The data are then sent to the data storage layer (cloud),

which contains scenario library, material library, antenna library, user database, task database, and measurement database. On this basis, the data analysis layer can perform channel simulation and modeling for the specified site using acceleration algorithms (e.g., parallel computing and space partitioning), thereby making the RT simulation more accurate and efficient. The output of the data analysis layer can then be applied to the system-level simulation, network planning, and antenna design. Thus, with the aid of our CloudRT platform, designers can accurately study the channel modeling in railway scenarios at both mmWave and THz frequency bands.

The electromagnetic calculation is composed of the different radio propagation mechanisms (i.e., LoS, reflection, scattering, transmission, and diffraction), which highly depend on the material of the object. Therefore, the material parameters are important to the accuracy of the RT simulation. The ITU provides dielectric parameters of typical building materials at 6–100 GHz [110]. Due to the complex composition of the dynamic scenario, it is unrealistic to measure the different propagation mechanisms of all materials. Therefore, the existing literature and works cannot cover all the material parameters that affect the propagation mechanism in all 5G-R scenarios. Given such a problem, Zhang *et al.* [111] proposed a method for calibrating the RT simulator based on channel measurements using a simulated annealing algorithm. The parameters of individual simulated MPCs are compared with those of the corresponding measured MPCs, and the ultrawideband electromagnetic parameters of typical materials in the target scenarios are obtained at the end. Based on the path loss measurement at 90 GHz in a railway viaduct scenario, similar calibration works are also performed in [112]. The difference between the fit path loss coefficients and the calibrated RT and the measurement is 0.01. With this, the system performance can be analyzed further for various configurations with higher accuracy. However, the accuracy of RT depends not only on the accurate material characteristics but also on the accurate geometry model of the propagation scenario and accurate superposition of multiple propagation mechanisms. The existing calibration methods assume that only a part of the measured propagation mechanisms exists and that the geometry of the environment model is built correctly. Besides, the referred measured MPCs mainly include the delay and power domain. In the future, researchers should consider all sorts of propagation effects and accordingly calibrate RT under different domains with the flexibility of adjusting the preconstructed geometry model. Owing to the development of high-resolution channel estimation algorithms, a sufficiently fine resolution of the measurement results is possible, thus allowing for a more accurate analysis of the impact of geometric features of the scenario and the material characteristics.

In terms of channel modeling, CloudRT supports the dual-mobility characteristics of a mobile scenario (i.e., dynamic scatterers and dynamic Tx/Rx). The high

computational complexity of RT has made research on its acceleration method a hot topic in the computer field. The typical approaches used include dimensionality reduction [113], space partitioning methods [114], ray-launching methods [115], and hardware-acceleration methods [116]. The basic principle of the dimensionality reduction method is to simplify the 3-D model in the elevation domain and then change the model to 2-D or 2.5-D. Although such method makes it much faster to compute multipath tracking, it makes the RT simulations less accurate. The idea in space partitioning is to divide the radio propagation scenario into several small regions such that researchers can avoid traversing all the objects in the scenario model for detecting MPCs. The space partitioning methods have no effect on the accuracy of the RT simulation; however, these methods mostly apply to a communication scenario where the scatterers are static. However, for 5G mobile communication scenarios for which the “dual mobility” feature exists, the appearance and disappearance of scatterers dynamically change the spatial distribution and the visibility relationship between the geometry planes, which makes performing RT more difficult. This also applies to ray launching methods. Accordingly, Nuckelt *et al.* [117] proposed a method for obtaining densely sampled channels by interpolation, which will lessen the RT simulation task and ensure that the simulation results are accurate. However, the interpolation algorithm in [117] focuses only on the amplitude domain and does not include the angular domain and polarization domain. Also, there is no theoretical analysis of reasonable sampling intervals. As the advantage of the hardware acceleration has already been proven in the aforementioned part, we need to combine new acceleration algorithms and hardware acceleration to make the RT calculation less complex.

Specifically, we propose to characterize the sparsity of dynamic scatterers and to determine the quasistationary interval in the specific scenarios. With regard to the sparsity characterization of dynamic scatterers, the global scattering center of the complex objects in the real mobile communication scenario can be extracted based on the full-wave analysis and image processing algorithms by setting a reasonable threshold. This method can be used to reveal the variation of the sparseness of the scatterer with frequency. Moreover, this method changes the modeling of complex objects (e.g., trains) from modeling deterministic polygons to multiparticle swarms that include explicit field strength and phase information in order to simulate RT accurately and efficiently under a complex dynamic scenario.

To determine the quasi-stationary interval in the spatial domain based on the geometric parameter transformation and spatial segmentation method, the CloudRT extracts the geometric features of a large number of similar scenarios from the scenario library. The relationship between the spatial geometric characteristics and the different propagation mechanisms is mapped through deep learning

methods. We can determine the quasi-stationary interval in the spatial domain using a predefined threshold and attain mathematical interpolation for different propagation mechanisms within a spatial quasi-stationary interval. Through this method, the sampling interval of dynamic RT simulations significantly lessens while maintaining accuracy. The railway channel characteristics at the 60-GHz band with 8-GHz bandwidth in the six scenario modules that use the two antenna height setups can be found in [118] and [119].

#### D. Massive MIMO Channel Modeling in Smart Rail Scenarios

Many research institutes and organizations have conducted channel measurements on massive MIMO under urban or suburban LoS and NLoS scenarios. The moving velocity of the train (Tx or Rx) is high in HSR scenarios. The virtual array technique for (massive) MIMO channel measurements cannot be applied to such time-varying environments (with the exception of using the train movement itself to create a massive virtual linear array at the train side). On the other hand, although one can apply a real array with demodulator chains to each receive antenna element, this would increase cost and would make it more difficult to implement. To this end, a switched array seems to be the best compromise between measurement time and hardware effort. In such a switched array, different antenna elements are connected to a demodulator chain (conventional channel sounder) via a fast RF switch. This would enable developers to measure dynamic massive MIMO channels. However, massive MIMO channel measurements in railway scenarios have the following five key challenges.

- 1) *Signal power or maintaining adequate SNRs and coverage:* The use of switched antenna arrays implies a lack of BF gain. The measurement SNR might thus be low, particularly in the cell edge of railway scenarios. A promising solution to address this problem is to use beam switching [120]. In this case, multiple antennas are activated simultaneously. If each antenna is equipped with a power amplifier, the effective conducted transmit power can be increased while realizing BF gain.
- 2) *Clock synchronization:* Clock synchronization is important in massive MIMO measurements to ensure that the frequency and/or signal phase of each channel is consistent with each other. In the dynamic railway scenarios, it is not possible to share the clock with the cable in the same way as that done in indoor measurements. Therefore, a standard GPS clock (to receive GPS pps signals from satellites) plus atomic clock, such as a rubidium clock, needs to be used in both Tx and Rx to keep precise synchronization.
- 3) *Sample rate:* In measuring dynamic mmWave channels under smart rail scenarios, the Doppler effect is expected to be significant due to high carrier

frequency and high movement speed. Thus, a high sample rate at the Rx is required to successfully measure the fast time-variant mmWave channel and acquire its Doppler characteristic. The Nyquist theorem in the time domain must be fulfilled. The sample for each subchannel of the mmWave massive MIMO system needs to be taken frequently to track changes in the dynamic channel. For a regular switching structure (each antenna measured in a regular sequence), this means that the switching cycle needs to be finished within the coherence time of the channel. This is challenging, particularly for mmWave massive MIMO channel measurements because the BS can have hundreds of antenna elements and more time is needed to measure all subchannels. Quasirandom switching sequences show promise to alleviate these problems; however, they have not yet been used in railway channel sounding applications. Moreover, the bandwidth of the excitation signal in the mmWave massive MIMO techniques always exceeds several hundreds of Megahertz. This large bandwidth also requires a high sample rate of the Rx.

- 4) *Fast switching:* As mentioned before, the measurements of the whole massive MIMO subchannels need to be finished within a limited time to ensure that all the measured data are within the same coherence time. Therefore, each subchannel needs to be quickly measured and fast switching is highly required. The measurement duration at one antenna element is determined by the actual subchannel measurement time and the time it takes to switch from one antenna element to the next. Obviously, the latter should be as short as possible, particularly in massive MIMO applications with hundreds of antenna elements. Thus, a switching system with a large number of ports, low insertion loss, and a short switching time is required. Note that since a single switch has limited RF ports, a cascade connection scheme of switches is required to build the above switching system. However, such scheme results in more challenging issues such as system calibration and synchronization. Moreover, Gb/s real-time data flow storage is needed in the measurement system to record data in real-time from the switching system.
- 5) *Dynamic BF tracking:* For mmWave massive MIMO measurements, BF technique is preferred to increase the gain. For the conventional static mmWave measurements with limited angular scanning range, horn antennas with large gains can be used for power concentrating. The transmitted signal beams are manually aimed at the Rx. In the dynamic scenarios of railways, how to accurately track the beam is significantly difficult. The channel sounding system should establish a feedback mechanism that can control the transmitter to adjust the beam angle according to the varying location of the Rx.

Based on the RT simulation platform, we mainly focus on analyzing the influence of massive MIMO configurations at the mmWave band inside the railway viaduct. The configuration is suggested by 3GPP and the propagation parameters are calibrated in our previous work [121].

### E. 5G Operating Frequency Bands for Smart Railways

Before we discuss potential frequency bands for the smart rail, we first need to reference the suitability of spectrum bands for the 5G generic services as defined by the European Union METIS: xMBB, mMTC, and ultrareliable MTC (uMTC).

- 1) xMBB requires a combination of frequencies, comprising lower bands below 6 GHz for both coverage and data traffic purposes and higher bands above 6 GHz with a large contiguous bandwidth. These requirements are needed to cope with the ever-increasing traffic demand, including WB solutions. An exclusive licensed spectrum essentially guarantees coverage and QoS, and this can be supplemented by a spectrum that other licensing regimes authorize to increase the overall spectrum availability. Such licensing regimes include licensed shared access or unlicensed access (e.g., WiFi offload) or new enhanced unlicensed access schemes (e.g., license-assisted access).
- 2) Some mMTC applications are most suited for frequency spectra below 6 GHz, particularly spectrum below 1 GHz, in cases when large coverage areas and good penetration are needed. Although an exclusive license is preferred, other licensing regimes may have to be considered depending on the specific application requirements.
- 3) Licensed spectrum is considered most appropriate for uMTC. For safety and for V2V and V2X communication, an option could be the frequency band 5875–5925 MHz harmonized for ITS. Another option could be the sub-1-GHz spectrum, which is particularly well suited for high-speed applications and rural environments.

NGMN requirements for 5G include the capacity to support a wide range of applications that may have differing requirements for the underlying mobile connectivity. This will consequently require access to a range of spectrum bands with differing characteristics in order to address a wide range of requirements for coverage, throughputs, and latency in the most cost-efficient manner and to enable the spectrum to be used effectively. Spectrum bands already licensed to mobile network operators will form an essential foundation for 5G mobile services. As such, it is important to allow operators to “refarm” existing spectrum bands to 5G technology according to their deployment strategy, which helps to improve SE and to introduce new capabilities. It will also help in planning the necessary long-term investments.

Further spectrum of MBB for both coverage and capacity may need to be organized in the future. A spectrum below 1 GHz is particularly useful for coverage, especially in indoor settings and rural areas. A spectrum above 6 GHz is particularly useful for supporting very high data rates and short-range connectivity.

In the view of the future demand of smart rail, the single bandwidth resource at 450 MHz will hardly be sufficient for the diverse potential services envisioned in evolved railway communication system. The available spectrum resources in TDD or FDD modes (at 20 or 25 MHz) can meet the concurrent traffic requirements of train-related services in all scenarios under the present communication system. In special scenarios (e.g., grand passenger train stations and marshaling stations), the frequency planning that supports massive-bandwidth requirement services adopts a spot-line combined approach. As such, this allows operators to attend to the application for shared licensed spectrum at 1.8 GHz for vertical customers or to utilize unlicensed spectrum at 2.4 or 5.8 GHz for complementation straightforwardly.

Researchers need to evaluate the possible candidate bands in higher frequencies to address the requirements for new spectrum for UDNs in 2020 and beyond. The industry needs such frequencies such that very wide bandwidth channels can support very high data rates and short-range mobile connectivity (e.g., 500–1000 MHz of contiguous spectrum per network to support the multitude of services). The total spectrum requirements also have to consider the potential need to accommodate multiple networks. Therefore, it is necessary to further develop the technical feasibility of the ranges between 6 and 100 GHz, particularly those ranges where primary allocations to mobile, as stated in the ITU radio regulations, already exist. Therefore, the lower limit for the band range (above 6 GHz) should be further assessed.

According to our simulation analysis, the low frequency band of 3.5 GHz and the mmWave band of 25–30 GHz can cover a 1-m railway line using directional antennas and reasonable transmitting power. The mmWave band can support higher transmission rate and larger bandwidth required by the dedicated communication for the smart rail scenarios. Considering spectrum utilization, smart rail services, propagation characteristics, BS deployment, frequent handover, and other factors, a reference frequency band 26–38 GHz can be adopted for smart rail communications.

## VI. 5G KEY TECHNOLOGIES FOR SMART RAIL

One of the key challenges for smart rail is achieving reliable high data rate transmissions at high speed under various railway scenarios. Providing transmission rates of up to Gb/s, and/or reliability of up to 99.9999% in high mobility is a very significant technical challenge, which is compounded in railway communications where speeds of up to 500 km/h occur. The physical-layer techniques

needed to achieve those goals are parallel to those that are used in 5G, in particular MIMO (including massive MIMO and distributed MIMO), as well as relaying and design of suitable codes and link margins for enhancing reliability. However, the special aspects of high mobility need to be considered. Sections VI-A–VI-H discuss all these aspects.

### A. ST-ASM in Railway Scenarios

SM is a digital modulation concept for multiantenna wireless systems, which has been introduced to increase SE while allowing for a low-complexity implementation. SM achieves multiplexing gain by mapping a block of information bits onto two information-carrying units: one symbol is chosen from a conventional constellation diagram and another unique transmit antenna index is chosen from the so-called spatial constellation diagram. At last, the point of the signal constellation diagram is transmitted through a single active antenna belonging to the spatial constellation [122]. At the Rx, the optimum maximum likelihood (ML) detector can simultaneously estimate the active Tx antenna index and the selected modulation symbol index. As a result, the overall SE increases (compared to transmission from a single antenna) by exploiting the transmit antenna index. Using SM in smart railway can ensure system robustness and reliable services.

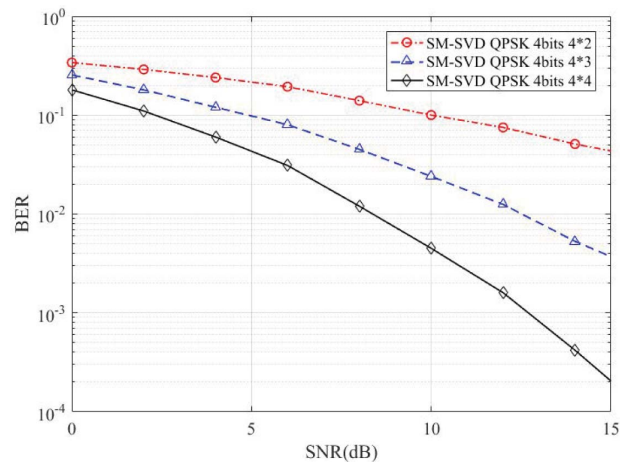
It follows from the above description that in SM, the data transmission rate is proportional to  $\log_2(N_t M)$  bits/s/Hz, where  $N_t$  is the number of the transmit antennas, and  $M$  denotes the MO. We can thus improve SE by increasing the number of transmit antennas. SM is different from a conventional MIMO transmission; in the latter, we can improve reliability through diversity and increase channel capacity by spatial multiplexing. In contrast, the former maps information to the Tx antenna index to obtain additional multiplexing gain.

When using SM, there is only one or a few active transmit antennas. Thus, only a single-RF chain is required; this reduces complexity and energy consumption, which make the commercial communication system more economical than conventional MIMO methods. This single RF scheme can also completely remove the intersystem interference at the Rx because only one Tx antenna is activated for any symbol duration. On the other hand, SM with multiple Rx antennas shows improved robustness because it offers excess degrees of freedom that can be used to cancel signals from man-made interference including intentional jamming.

The transmission of an  $N_t \times N_r$  MIMO system can be represented as

$$y = \sqrt{\rho} H x_q + w \quad (1)$$

where  $\rho$  is the average SNR at the transmit antenna and  $H$  represents the  $N_r \times N_t$  channel matrix, in which  $h_l$  is the  $j$ th column vector of the channel matrix. Finally,  $w$  denotes the  $N_r$ -dimension AWGN with i.i.d. elements



**Fig. 8.** BER performance versus SNR evaluated using SM-SVD approach for the case of a 4-b/s/Hz transmission with different receive antennas. The number of the transmit antenna is 4 and the number of the receive antenna is 2, 3, and 4, respectively. All of the three schemes employ a QPSK modulation alphabet.

and  $x_q$  denotes the transmitted symbol. The decision of the ML detection is given by

$$(\hat{J}, \hat{q}) = \arg \min_{j,q} \|y - h_j x_q\|_F^2. \quad (2)$$

Considering the large additional gains of SM with more Tx antennas, railway communications might be an especially beneficial application of SM. The roof of the train carriages can provide enough space; to deploy a few hundred antennas. Li et al. [123] proposed a different detection method called the SM-SVD. This method has a low complexity and is suitable for a massive MIMO system. The detection method includes two independent estimation steps and the estimation of the antenna index can be written as

$$\hat{J} = \arg \min_j \arccos \frac{|\langle h_j, y \rangle|}{\|h_j\|_F \|y\|_F} \quad (3)$$

where  $\langle \bullet \rangle$  denotes the inner product in the Hilbert space and  $\|\bullet\|_F$  represents the Euclidean norm. Once the index of the transmit antenna is estimated, the transmitted symbol can be estimated in the traditional way.

Let us assume the channels are independent Rayleigh flat-fading and the noise is additive Gaussian. We then evaluate the SM using different receive antennas for a 4-b/s/Hz transmission. Fig. 8 shows the numerical results when we employ the SM-SVD. We can see that the SM scheme with four Tx antennas performs the best, which demonstrates that the extra Rx antennas can compensate for the loss of SNR, and thus bring additional gains [123].

Most of the recent works on SM have focused on generalized SM and ASM. Unlike SM, in which only one antenna is used in each time slot, two antennas are activated

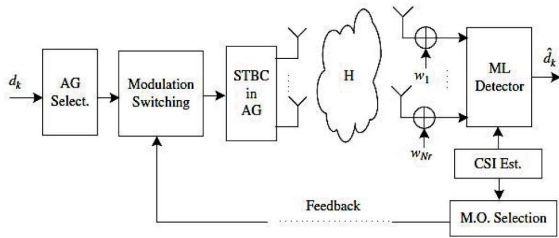


Fig. 9. Block diagram of ST-ASM transceiver.

simultaneously in generalized SM to transmit information symbols, thus transmitting multiple data streams. The data rate of a generalized SM system is larger than that of a conventional SM assuming an equal number of available Tx antennas. In ASM, different levels of modulation are selected adaptively based on a low-complexity MO selection criterion for different channel conditions. Thus, we achieve better BER performance while using the same data rate target when using ASM than using conventional SM methods with fixed MO.

To be robust in smart railway scenarios, an ASM scheme was proposed in [124], in which the authors embedded an Alamouti STBC to the high-speed communication system. As shown in Fig. 9, in the ST-ASM system, two antennas are activated to transmit Alamouti code blocks, which are obtained through encoding information symbols. Compared with the conventional SM and ASM, the proposed ST-ASM scheme provides larger diversity gain and better BER performance when using equal data rate and equal transmission power [124]. Similar to the ASM scheme, the MOs assigned to each Tx antenna under the ST-ASM are chosen adaptively according to different channel conditions. Note that there is zero-order modulation present in the antenna group when selecting MOs and no information bits are modulated at the Tx antennas. As such, we consider the ML detection of the ST-ASM.

## B. Fast Time-Varying Channel Estimation in HSR Scenarios

HSR trains run at a speed of around 360 km/h, which leads to rapid changes in both amplitudes and phases of wireless channels. As a result, fast time-varying channel estimators are of vital importance to the wireless transmission systems of HSRs. Classical channel estimators (i.e., LMMSE and LS) may fail to provide satisfactory performance in such highly time-selective scenarios.

Traditionally, there are two widespread ways to estimate the time-varying channel parameters: through a BEM and through a GMM. The former decomposes the channel parameters into a superposition of the time-varying basis functions, which is weighted by time variant coefficients; and the latter tracks channel variation through symbol-to-symbol updating.

BEM expresses the time-varying channel  $\mathbf{h} = [h_1, \dots, h_N]^T$  as  $\mathbf{h} = \mathbf{B}\mathbf{g}$ , where  $\mathbf{B} = [b_1, \dots, b_Q]$  is the basis matrix,  $b_q (1 \leq q \leq Q)$  denotes the  $q$ th basis vector, and  $\mathbf{g} = [g_1, g_2, \dots, g_Q]^T$  contains the BEM coefficients to be estimated. Clearly, the number of the channel parameters to be estimated is decreased from  $N$  to  $Q$ . Often choices of BEMs are complex-exponential BEM and polynomial BEM. However, these BEMs generally have unvaried basis matrices, which may not be optimal for various HSR scenarios such as cuttings, viaducts, and tunnels. Similarly, GMMs use the same model for all HSR scenarios.

It is worth noting that trains run on a fixed track, which is different from mobile users on the ground that can move toward any directions. Therefore, the communication systems of the current train face the same wireless environment with those of the past trains, which can be further exploited to aid channel estimation. Because of this, we propose applying a new channel estimator for HSR that exploits the historical information of past trains to enhance the estimation performance [80], [125]. This channel estimator is referred to as HiBEM.

The underlying principle of HiBEM is that trains follow fixed tracks and that the surrounding geological and wireless environments are almost the same at a fixed position given similar weather conditions. Consequently, the wireless channels of the current and past trains in certain areas are strongly correlated if the transceiver configuration is known and that train speeds are predictable. Based on these observations, HiBEM adopts the first  $Q$  eigenvectors of the channel autocorrelation matrix as its basis matrix. The channel autocorrelation matrix is defined as  $R_h = E\{hh^H\}$ . The eigenvalue decomposition of  $R_h$  will yield  $R_h = UDU^H$ .

Hence, we can construct basis matrix  $B$  from  $U$ , which has its first  $Q$  columns as the eigenvectors of  $R_h$ . The basis matrix  $\mathbf{B}$  of HiBEM is chosen as  $\mathbf{B} = U(:, 1:Q)$ . It has been proven in [148, Sec. III-B] that basis matrix  $\mathbf{B}$  is optimal in terms of minimizing the approximation MSE, that is,  $\min E\|h - B\eta\|^2$ , where  $E$  denotes average operation and  $\eta$  is a coefficient vector.

One key problem is how to compute the channel correlation matrix  $R_h$ . The wireless channels on HSR are strongly related given a fixed position, which can be used to calculate the correlation matrix  $R_h$ . Fig. 10 depicts the scheme of calculating  $R_h$ . Suppose the current train is the  $k$ th train. We utilize the estimated channels of previous  $(k-1)$  trains at the same position to compute the corresponding correlation matrices  $R_{hi} (1 \leq i \leq k-1)$  and average them as

$$\mathbf{R} = E\{\mathbf{R}_{hi}\} = \frac{1}{k-1} \sum_{i=1}^{k-1} \mathbf{R}_{hi}. \quad (4)$$

We utilize  $\hat{R}_h$  to derive the real channel correlation matrix  $R_h$ .

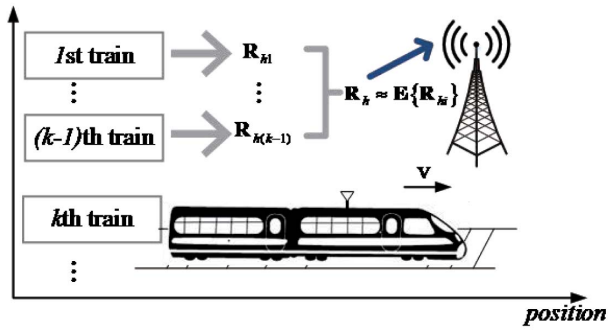


Fig. 10. Scheme of estimation with historical information.

Fig. 11 compares the estimation performance of various BEMs. We set the channel vector length as 20 and generate random channels through a Jakes model with a maximum Doppler shift of 300 Hz. Then we choose CE-BEM and HiBEM to obtain the channel estimates and their corresponding MSE. For comparison, the estimation MSE of the interpolation method is also plotted. It can be readily checked in Fig. 11 that HiBEM has the minimum approximation error and the error decreases with larger  $Q$  value.

Fig. 12 depicts the estimation performance of HiBEM, LS, and LMMSE approaches in terms of their mean square errors; parameter  $k$  indicates the number of past trains. We consider an OFDM system with 64 subcarriers, that is,  $N = 64$ . We set the number of pilots in the LS and LMMSE estimators to 16. The train moves at 360 km/h, which indicates that the train moves 0.1 m in 1 ms or during one LTE frame duration. The communication system is deployed at 3.5-GHz carrier frequency with 100-MHz bandwidth. The transmitted symbols are generated through BPSK modulation. The white Gaussian noise variance is set

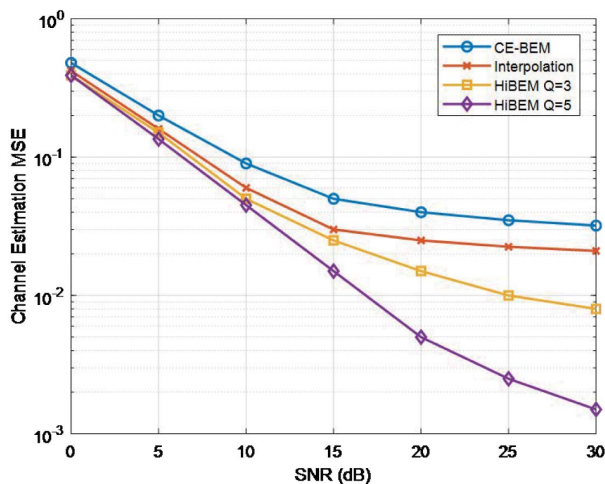


Fig. 11. Performance comparison of various BEMs.

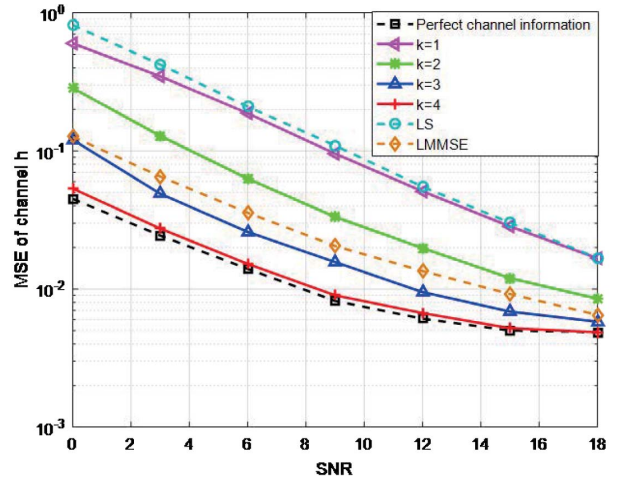


Fig. 12. Estimation performance of HiBEM, LS, and LMMSE approaches.

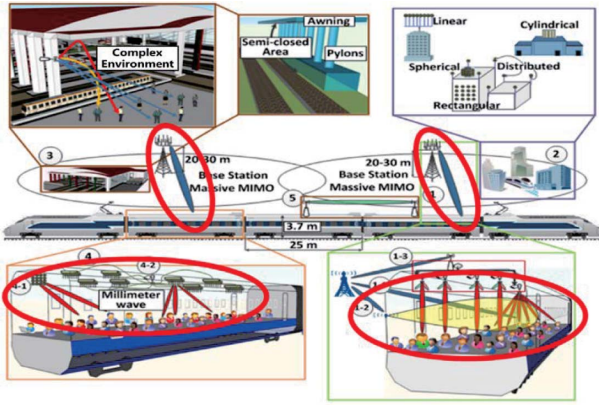
to 1. It can be shown from Fig. 12 that HiBEM outperforms both LS and LMMSE. It can also be seen that HiBEM, in the case of four trains, achieves almost the same estimation performance compared to the case when the channel correlation matrix is perfectly known at the Rx.

### C. Massive MIMO for HSR

When a train arrives in a slow-moving or stationary “in-station communication” scenario, we can upload the temporarily stored data to the train using large-scale mmWave arrays. Accordingly, we can obtain the traffic data, passenger information data, and multimedia data required for the next travel section. As for the hotspot railway areas (e.g., railway marshaling stations and hub areas), we can use other broadband technologies as a supplement. The users in these areas move slowly; thus, we can use high frequency and broadband modulation technologies to provide large-capacity wireless communication services.

Massive MIMO extends the multiuser MIMO concept by dramatically increasing the number of antennas employed at the BS to be able to serve more users simultaneously within the same time-frequency block. Note that massive MIMO does not significantly increase the peak rate for a single user; it inherently obtains its high SE by serving multiple users simultaneously. Fig. 13 illustrates the applications of mmWave massive MIMO under in-station and in-carriage scenarios.

With a large number of antenna units, the centralized massive MIMO BSs of railway systems can use BF to form high-gain and directional narrow-beam targets to track high-speed moving trains, as shown in Fig. 14. By using a high-gain antenna, the SNR and channel capacity improves. Multiuser spatial division multiple access can be achieved by using narrow beams, which will accordingly

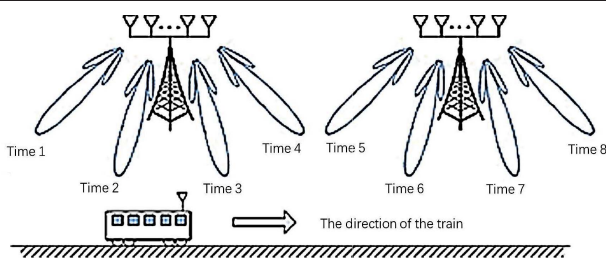


**Fig. 13.** mmWave massive MIMO under in-station and in-carriage scenarios.

increase the overall capacity of the BSs. Beam-tracking is utilized to achieve dynamic focusing of the energy to the intended Rx. The distributed massive MIMO configuration inside the carriages can provide high data rate services for multiusers.

The feasibility of using massive MIMO in several HSR scenarios has been investigated in [126]. The key challenges to its practical implementation include the requirement for fast signal processing in light of the high train mobility and the low EE of concentrated large-scale arrays along HSR lines. One of the promising solutions to tackle these key challenges is to utilize a distributed massive MIMO architecture called cell-free massive MIMO [127].

Unlike the conventional collocated massive MIMO, cell-free massive MIMO employs a large number of distributed APs with antennas over a wide area. First, the required net throughput and QoS are quite high inside a very large railway station since the station may contain thousands of passengers. The goal then is to provide the same high data rate to all passengers while using limited transmit power. Second, along the railway lines, cell-free massive MIMO can be deployed to guarantee good coverage without cells or cell edges. Remember that in conventional cell-based systems, the receive signal power drops to a minimum when the UE is near the cell edge, which may lead



**Fig. 14.** Block diagram of the beam tracking under railway scenarios.

to low service quality at those locations. However, the QoS of cell-free massive MIMO is uniformly good for all areas. Finally, the frequent handover problem (e.g., hundreds of UEs need to be handed over to the next BS in a very short time) can be eliminated in a wide area.

We use  $q_k$  to denote the information symbol of the  $k$ th UE and  $\rho_u$  to denote the maximum transmit power. The channel response  $g_{mk}$  between AP  $m$  and UE  $k$  is assumed to be Rayleigh fading as  $g_{mk} \sim \mathcal{CN}(0, \beta_{mk})$ . With the assumption of perfect transceiver hardware, the received signal sequence at the  $m$ th AP is

$$y_{um} = \sum_{k=1}^K g_{mk} \sqrt{\rho_u \gamma_k} q_k + w_{um} \quad (5)$$

where  $\gamma_k$  denotes the uplink power control coefficient and  $w_{um}$  represents the additive noise. We introduce  $\kappa_r$  and  $\kappa_t$  to denote the level of hardware impairments in the Rx and transmitter, respectively. Then, the uplink SE of the  $k$ th UE can be given by

$$R_{uk} = \log_2 \left( 1 + \frac{\kappa_r \kappa_t A_u}{\kappa_r B_u + \kappa_r C_u - \kappa_r \kappa_t A_u + (1 - \kappa_r) D_u + E_u} \right) \quad (6)$$

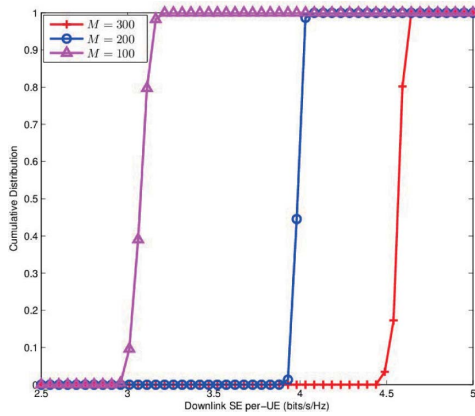
where

$$\begin{aligned} A_u &\triangleq \gamma_k \left( \sum_{m=1}^M \lambda_{mk} \right)^2 \\ B_u &\triangleq \sum_{k'=1}^K \gamma_{k'} \left( \sum_{m=1}^M \lambda_{mk} \beta_{mk'} + \rho_p (1 - \kappa_r) \sum_{m=1}^M c_{mk}^2 \beta_{mk'}^2 \right) \\ C_u &\triangleq \sum_{k'=1}^K \gamma_{k'} \left( \left| \varphi_k^H \varphi_{k'} \right|^2 + \frac{1 - \kappa_t}{\kappa_t \tau} \right) \left( \sum_{m=1}^M \lambda_{mk} \frac{\beta_{mk'}}{\beta_{mk}} \right)^2 \\ D_u &\triangleq \sum_{m=1}^M \left( \lambda_{mk} \sum_{k'=1}^K \gamma_{k'} \beta_{mk'} + c_{mk}^2 (1 - \kappa_r) \rho_p \beta_{mk'}^2 \right. \\ &\quad \left. + c_{mk}^2 \kappa_r \rho_p \beta_{mk'} \left( \tau \kappa_t \left| \varphi_k^H \varphi_{k'} \right|^2 + (1 - \kappa_t) \right) \right) \\ E_u &\triangleq \frac{\sigma^2}{\rho_u} \sum_{m=1}^M \lambda_{mk}. \end{aligned} \quad (7)$$

In the following, we numerically investigate the SE of cell-free massive MIMO systems with hardware impairments. We assume  $M$  APs and  $K$  UEs are independently and uniformly distributed in an area of size  $1 \times 1$  km<sup>2</sup>. We further assume that the number of pilot sequences is equal to the number of UEs, and all UEs are assigned with orthogonal pilots. Therefore, the pilot contamination can be alleviated. The detailed simulation parameters can be found in [128].

Fig. 15 shows the cumulative distribution of the downlink SE per UE for cell-free massive MIMO for the different numbers of APs (e.g.,  $M = 100, 200, 300$ ). We assume





**Fig. 15.** Cumulative distribution of the downlink SE per UE for cell-free massive MIMO against different number of APs.

that the number of UEs is 10. Clearly, the SE performance improves significantly in terms of both median and 95%-likely per-user SE for a larger number of APs. It is clear that the power control algorithm proposed in [128] can significantly improve the SE performance of cell-free massive MIMO systems.

However, applying cell-free massive MIMO to HSR scenarios also poses the following major challenges.

- 1) The fronthaul requirements, in terms of both capacity and power, increase significantly. One should be careful in designing signal processing techniques that reduce data traffic between APs and the central processing unit. The tradeoff between the backhaul requirements and the system performance for cell-free massive MIMO is noteworthy, and thus needs to be further studied in the future.
- 2) The coherence times of the channel impulse responses in high-mobility scenarios are significantly small. The deep learning-based channel estimation schemes proposed in [129] and [130] can achieve remarkable performance. Thus, the current number of pilots is not enough to enable orthogonal pilots from all APs; yet reusing the same pilots may cause pilot contamination [131]. Therefore, the effect of pilot contamination on future 5G-R needs to be further investigated.

Cell-free massive MIMO and MCBF are two promising technologies for HSR scenarios. However, due to their characters, they can be employed in different conditions. For example, MCBF utilizes several proximate BSs to coordinately serve the UEs onboard high-speed train. There still exists frequent handover between different BS clusters. In contrast, cell-free massive MIMO eliminates handover by using efficient scheduling and high-capacity fronthaul. The complexity of cell-free massive MIMO is more pronounced than the one of MCBF, while achieving better performance. MCBF can be used for low-speed scenarios, such as rolling in and out of a railway station.

#### D. mmWave Massive MIMO in Smart Rail Scenarios

mmWave communications employ massive MIMO mainly to achieve directional transmission based on antenna array BF. Accordingly, a major topic that researchers focus on is how to reduce the implementation complexity of mmWave railway communications without significantly degrading its performance. Based on the available velocity and location information that the communication-based train control system provides, we present a new simplified beam-tracking scheme that ensures high-capacity and low-outage probability. Unlike the beam-switching methods standardized in IEEE 802.11ad or the fixed selection of beamwidth as a function of location uncertainty, we can develop and optimize the sum-rate capacity maximization problem by jointly adjusting the beam direction and the beamwidth, and then introduce high flexibility in the beam patterns. Furthermore, we can apply a genetic algorithm to solve the decomposed problems and to obtain the near-optimal beamwidth in order to determine the beam direction [132]. In summary, the beam direction acts as the leader with a predefined azimuth angle, and we aim to maximize the sum-rate capacity. Thus, the joint beam direction and width allocation problem can be treated as a leader-follower sequential decision problem. The beam direction will not change, unless the transmit antenna gain cannot be improved and QoS requirements will not be satisfied by a single beam with width control. When the optimal beamwidth appears to be far from the QoS requirements, a dynamic beam direction adjustment is sought to cope with the transceiver angle, and the BS to MRN link turns the centerline of the main lobe. Thus, the implementation complexity is significantly reduced, while the sum rate remains close to optimal.

We consider a typical system model of BF operating at mmWave for railway wireless communication systems as depicted in Fig. 16. The figure shows that the trackside BS is equipped with a ULA array. A moving relay node is also mounted on top of the train to help enhance the cellular coverage and to avoid vehicle penetration loss [132]. Then, we present the performance evaluation results of our beam-tracking scheme. In the simulations, a typical mmWave small cell in the HSR scenario is considered. Special propagation properties, high mobility, and the Doppler effects are also considered.

For comparison, two typical BF strategies, real-time beam tracking and beam switching, are conducted. Specially, the real-time beam tracking strategy is configured with the beamwidth setting as  $10^\circ$  and  $25^\circ$  separately, wherein the beam direction is always aligned with the MRN when the QoS requirements cannot be satisfied. Besides, the beam switching scheme in IEEE 802.11ad is also considered for comparison, where  $N_b = 3$  candidate beam patterns are predefined with equal angle interval and fixed beamwidth is set as  $10^\circ$  and  $25^\circ$ . In the simu-

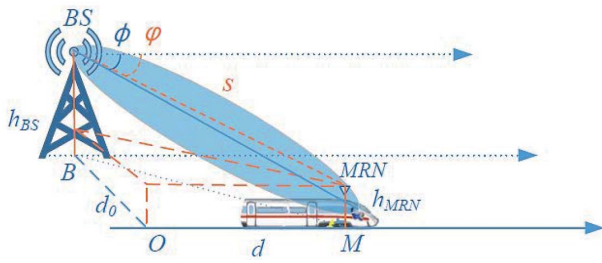


Fig. 16. mmWave massive MIMO beam-tracking systems for HSR communications.

Table 3 Simulation Parameters for Massive MIMO BF

Parameter	Symbol	Value
System bandwidth	W	2160 MHz
Background noise	$N_0$	-174 dBm/Hz
Path loss exponent	$\alpha$	2
Carrier frequency	$f_c$	60 GHz
Height of the BS	$h_{BS}$	4 m
Height of the MRN	$h_{MRN}$	2.5 m
Perpendicular distance	$d_0$	5 m
Radius	$R$	580 m
Threshold	$\gamma_{th}$	10 dB

lations, the proposed strategy is configured with a valid range of a beamwidth from  $10^\circ$  to  $25^\circ$ . The simulation results are depicted in the following four figures, which demonstrate the sum rate and outage of the proposed algorithm under various configurations, including transmission power, velocity of the train, and the train-to-BS distance. The simulation parameters are given in Table 3.

Fig. 17 plots the received SNR from different train positions under three BF schemes. The received SNR decreases as the train moves away from the BS. The rate of decrease is higher when the train is near the BS. When the train moves farther away from the BS, on the other hand, the SNR decreases more slowly. The SNR of beam switching varies significantly with train position since the relative position of the BS and the relay on the train decide the beam direction. The SNRs of the real-time beam tracking and our proposed BF scheme vary less since the beam direction adjusts according to the relative position of the BS and the relay.

Fig. 18 plots the sum rate of the proposed scheme under different train velocities and BS transmission powers.

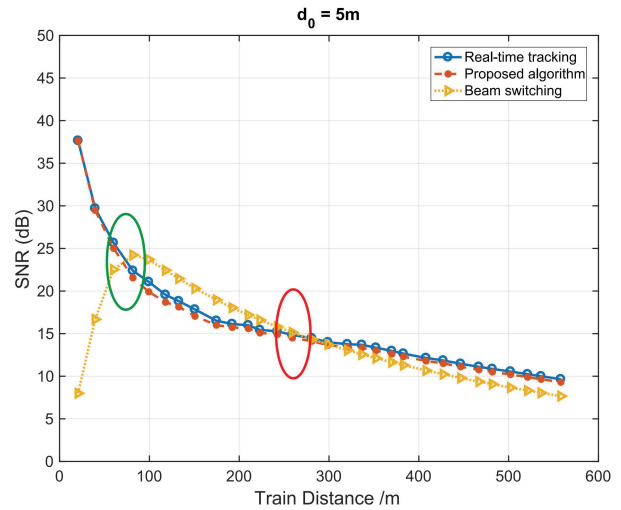


Fig. 17. Received SNR from different train positions.

The sum rate is the total data rate between the BS and the relay in the serving cell. Fig. 18 shows that as the velocity of the trains increases, the Doppler effect consequently becomes more serious, and the sum rate decreases. When the transmission power increases, the received SNR increases; thus, the sum rate also increases.

Fig. 19 depicts the sum rate comparison of the three schemes versus train velocity with  $P_T = 30$  dBm. We can observe that the sum rate drops quickly when the train speed increases, owing to the high ICI introduced by the more severe Doppler spread.

Fig. 20 displays the outage performance of these BF techniques versus train velocity for  $P_T = 10$  dBm. The outage probability rises with the speed of the train, owing to aggravating ICI introduced by Doppler spread. The gap between the proposed algorithm and the real-time tracking is indistinguishable at the speed of 100 km/h, and rises to

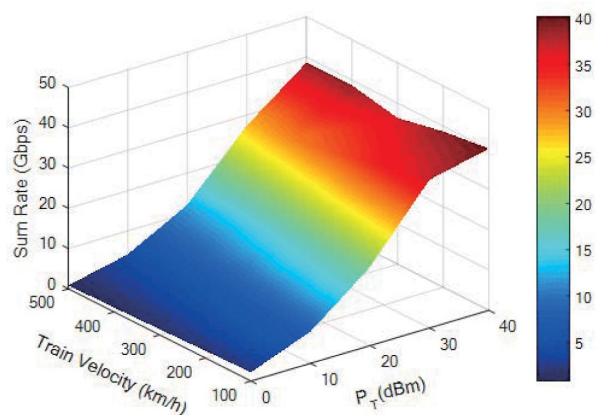


Fig. 18. Sum rate under different train velocities and BS transmission power.

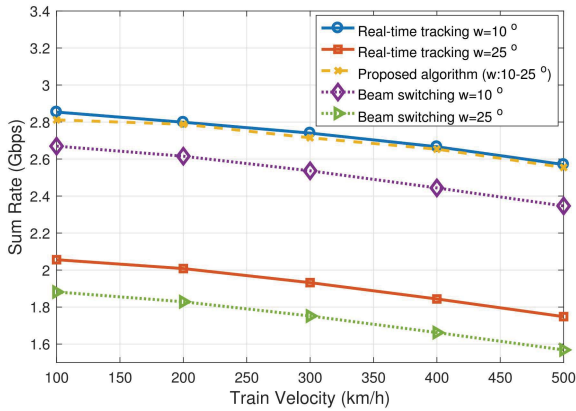


Fig. 19. Sum-rate comparison of BF strategies at different v.

8.69% at a speed of 500 km/h. However, the counterpart gaps between the proposed one and the beam switching are 73.92% and 47.23%, respectively.

Considering varying packet loss due to the high speed of HSRs and possible mmWave obstacles, the performance of TCP-based applications is also important. It is very interesting to see how common user applications that use TCP might be affected with the current architecture for future smart railways. This will be done in our future work.

### E. Efficient BF for High-Speed Train in Low-Mobility User Existence Scenario

Hybrid BF is a technique in which a large number of antenna elements are connected, via an analog BF matrix in the RF domain, to a smaller number of RF chains that connect to the BBUs [133]. Thus, data transmission can be achieved by a combination of digital BF and analog BF. This technique can drastically reduce cost and energy consumption as compared to employing one RF chain for each antenna element, while resulting in only minor performance loss under most circumstances. It is useful at

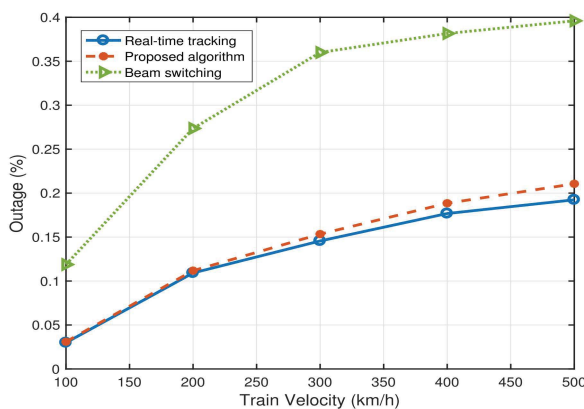


Fig. 20. Outage analysis of BF strategies at different train speed v.

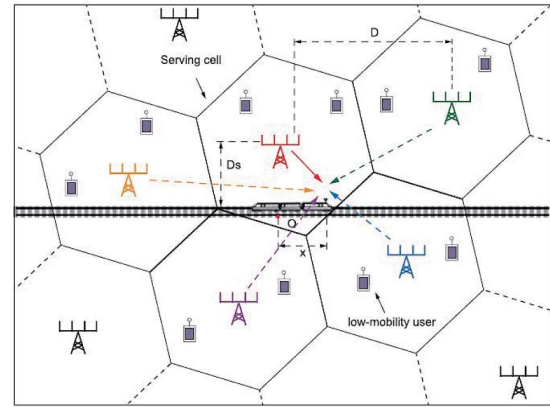


Fig. 21. HST goes through the mobile cellular system with low-mobility users.

all frequency ranges but particularly important at mmWave frequencies [134].

We now design a robust and efficient communication scheme with mmWave hybrid BF. First, we establish a movement model for railways. Thereafter, we analyze the change in channel state when the train is moving, and then we design the channel-state monitoring mechanism, the beam tracking mechanism, and spatial diversity. We evaluate the performance of the proposed BF/tracking mechanism and compare it with the system performance without it.

The key issue for future mmWave studies is “How can we ensure a robust and efficient data transmission between BS and the train carriage mmWave links while considering the movement of the train carriage?” Specifically, this issue includes answering the following questions: How can we design an efficient hybrid BF to achieve multistream parallel and efficient transmission of hybrid BF technology? How can we improve the robustness of the connection using spatial diversity technology? How can we adapt the corresponding technologies for the hybrid BF system based on mobile features, such as analog BF, digital precoding, and spatial diversity combining?

Fig. 21 depicts a typical scenario in which a high-mobility train passes through a mobile cellular system that also serves some traditional low-mobility users outside the train. That is to say, the high-mobility train and low-mobility users coexist in a single system. To improve SE, we employ multiantenna BSs to serve the high-mobility train and low-mobility users over the same frequency band at the same time. Therefore, the interuser interference between the train (i.e., the big user) and low-mobility users and the intercell interference among different cells cannot be neglected in HSR communication systems [135].

MCBF is an effective way for mitigating both intercell and intracell interference since it is able to exploit the spatial degree of freedom to concentrate the beamformers to the users of interest. Therefore, MCBF may also be an efficient way of improving the information-transmit

performance of railway communications with high mobility. Generally, existing MCBF designs are classified into two types: the power minimization-oriented design and the max–min fairness-oriented design. However, both of them may no longer be efficient in HSR communications. First, the primary goal of 5G-R is to increase the data rate for the train passengers rather than to save energy. Besides, BSs along the railway line often have stable and sufficient power supply. Thus, minimizing power is not the main concern of HSR communications. Second, to avoid group handover, all train passengers are aggregated to form a single virtual high-mobility user (i.e., the big user), which requires much higher data rates than low-mobility users have. Thus, in this situation, a max–min fairness between the big user and the low-mobility users is not meaningful. Therefore, existing MCBF designs are not suitable for HSR communications.

Consequently, we propose an improved MCBF design for HSR communication systems that is capable of providing the train with the highest information rate possible. The considered design is formulated into an optimization problem with the objective to maximize the achievable rate of the train under the constraints of minimal rate requirement of low-mobility users and the power budgets of all BSs

$$\begin{aligned}
 & \max R_{s,0}(w_{s,0}, \dots, w_{N_C,K}) \\
 & \text{s.t. } R_{s,K}(w_{s,0}, \dots, w_{N_C,K}) \geq r_{s,k}, \quad k = 1, \dots, K \\
 & \quad R_{n,K}(w_{s,0}, \dots, w_{N_C,K}) \geq r_{n,k}, \quad k = 1, \dots, K \\
 & \quad \sum_{l=0}^K \|w_{s,l}\|_2^2 \leq P_T \\
 & \quad \sum_{l=0}^K \|w_{m,l}\|_2^2 \leq P_T, \quad n = 1, \dots, N_C
 \end{aligned} \tag{8}$$

where  $s$  and  $n = 1, \dots, N_C$  are the index of the serving BS and the  $n$ th neighboring BS, respectively. In each cell,  $k$  is used to indicate the  $k$ th user, where  $k = 0$  indicates the train and  $k = 1, \dots, K$  indicates the low-mobility user. Here,  $w_{x,y}$  denotes the BF vector from the  $x$ th BS to its  $y$ th user and  $r_{x,y}$  denotes the minimal rate requirement of  $y$ th low-mobility user in the  $x$ th cell. Finally,  $P_T$  denotes the power budget of each BS. For the high-speed train, the Doppler effect can be regarded as the source of ICI and the Doppler interference factor can be calculated by

$$P_{\text{ICI}} = 1 - \int_{-1}^1 (1 - |\tau|) J_0(2\pi f_d T_s \tau) d\tau \tag{9}$$

where  $f_d$  is the maximum Doppler frequency,  $T_s$  is the symbol duration, and  $J_0(\bullet)$  is the zeroth-order Bessel function of the first kind. Note that the proposed design involves the fairness among low-mobility users. As the goal of the proposed design is to maximize the achievable rate of the train, the rate requirement constraints on the low-mobility users shall hold with equality when the optimal

result is obtained. The reason is that the higher rate of the low-mobility users, the more interference goes to the train. By varying the minimal rate requirements of low-mobility users, a balance between fairness and performance can be achieved.

The considered problem is nonconvex, and an algorithm based on a bisection method is presented in [135] to globally solve it when the perfect CSI of the train and the low-mobility users is available at all BSs. The main idea of the solution approach is to use the bisection method to find the maximum data rate that the train can receive while all constraints can be satisfied.

In practice, due to the high mobility, the CSI of the train is hard to collect. Thus, the imperfect CSI of the train should be taken into consideration. To provide a high quality of wireless coverage, a broad BF approach is presented in [135] for the HSR communication systems. The algorithm presented in [135] can easily achieve a broad BF coverage for the train by taking the channel estimated error of the train into consideration. As the channel estimation error increases (which happens when the train moves with higher velocity), the BF coverage for the train becomes broader and the power consumption at the BSs increases. How to collect the CSI is also of high importance for transmit design, which is, however, beyond the scope of this article. The proposed design can be extended to the scenarios with other channel assumptions and the insights based on this article shall hold true.

The simulation scenario is given in Fig. 21. The number of neighboring cells is two, the number of transmit antennas is six, and the number of low-mobility users in each cell is two. The transmit power at each BS is set to 30 dBm. The required SINR at each low-mobility user is 5 dB. The system bandwidth is 10 MHz, and the noise power spectral density is  $-162$  dBm/Hz. The inter-BS distance  $D$  is 1000 m, and the distance between the BSs and the track is 300 m. For the train with 100 m/s, the relevant Doppler interference factor is estimated to be  $-19.4$  dB. The locations of two BSs are set as  $(0, 300)$  and  $(1000, 300)$ , respectively. We consider the case that the train moves from  $(-500, 0)$  to  $(500, 0)$ . For a given location of the train, the distance between each BS and the train can be calculated and the channel between each BS and the train can be estimated. Then, the optimal transmit BF design can be achieved by utilizing the algorithm presented in [135]. Some metrics are considered to evaluate the performance of the proposed design, such as mobile service denoting sum rate achieved at the train between two locations and the average information rate during this period.

Fig. 22 shows the mobile service amount of the train and the low-mobility user. It is observed that the derivative of the service amount of the train first increases and then decreases with regard to the train location whereas that of the low-mobility user does not change. Moreover, increasing the transmit power  $P_T$  improves only the service amount of the train. The reason is that our design intends

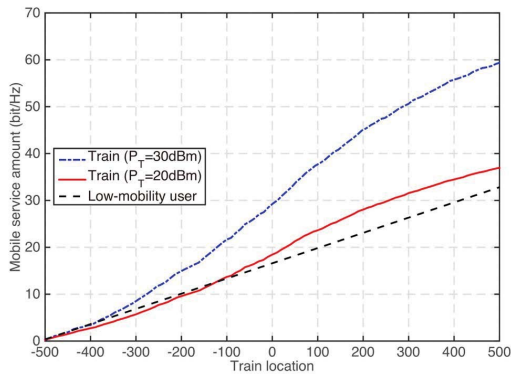


Fig. 22. Mobile service amount versus the position of the train.

to maximize the received rate of the train. Once the required minimal information rates of low-mobility users are satisfied, all remaining power of the serving BS should be allocated to increase the information rate of the train. Meanwhile, in order to mitigate the interference to the train, the remaining power of the neighboring BSs should not be used.

### F. MWB for Smart Rail Systems

In many cases, both rural and urban, an MWB is often the only technologically or economically feasible backhaul method. Even when an optical fiber is available near a cell site, MWB allows faster deployment until the fiber can be installed. On the other hand, it is expected that 5G will use spectrum above 6 GHz. Therefore, ideally, any spectrum used for 5G access should be flexible enough such that it can also be used for a backhaul. Likewise, spectrum licenses should be flexible enough to allow operators to meet the rollout demand while being capable of using 5G spectrum for backhaul whenever appropriate.

Higher frequency bands, especially millimeter and sub-millimeter wave bands, often show higher directionality, and thus can be used with more directional antennas, than sub-6-GHz channels. Hence, this makes it easier to use these bands simultaneously for mobile access and backhaul. Yet the specifics of the orthogonalization between these two types of transmissions still need to be further studied.

A related problem arises when the AP on the train aggregates the information from different users, and relays them to the BS. The latter part can be seen as a mobile backhaul from the AP to the BS; clearly in this case an MWB is the only option.

Fig. 23 shows an MIMO-BS only architecture and our proposed MWB architecture. The MIMO-BS is equipped with  $N_t$  antennas and is located on the track side, with  $N_u$  UEs located inside the carriage. In the MIMO-BS-only architecture, the MIMO-BS directly connects to the UEs in a multiuser MIMO mode using ZF-BF. In the MWB

architecture,  $N_s$  APs are deployed from above the carriage; the access link can be AP-UE and MIMO-BS-UE. The backhaul link between the AP and the MIMO-BS operates in a multiuser MIMO mode using ZF-BF. The MWB architecture uses a hierarchical joint user association and a backhaul bandwidth-allocation algorithm, which has an appropriate BF capacity model. We then compare the data rates of both architectures using the same parameter setups.

At the same time, the propagation between the  $N_t$  BS antennas and the  $N_s$  APs can be described through an  $N_t \times N_s$  matrix, in which for simplicity the entries of the small-scale fading are assumed to be independent  $\mathcal{CN}(0, 1)$  random variables. We further assume that the MIMO-BS antenna size  $N_t$  is much larger than  $N_s$  to keep our problem within the regime of massive MIMO. For convenience, the link between the BS and UEs denotes BS access link, the link between APs and UEs denotes AP access link, and the link between BS and APs denotes backhaul link. We define the set of UEs by  $U$  and the set of APs by  $S$ . Thus,  $S \cup 0$  is the set of BS and APs, in which the index 0 is introduced for the BS. We assume a ZF algorithm to reduce the interferences caused by inter-UEs and inter-APs. To fully utilize the advantages of massive MIMO, the MIMO-BS concurrently serves a backhaul link and a BS access link in a downlink time slot. The maximum size of the precoding group that the MIMO-BS can serve is  $N_g$  ( $N_g \ll N_t$ ), with  $N_g$  larger than  $N_s$  such that the backhaul links for all the APs can be accommodated over the same frequency band, where  $N_u$  is much greater than  $N_g$ . Resource sharing must be used when a large number of UEs are associated with the massive MIMO-BS. By allocating different orthogonal frequency-temporal resources to different precoding groups, which consist of maximum concurrent serving terminals, the system has no intergroup interference. After the APs establish the backhaul link, the data are delivered to the UEs from the AP access link. TDMA can be applied to this conventional point-to-point communication scenario.

We consider the mmWave bands 18, 28, and 38 GHz, each with 1-GHz bandwidth. The number of antenna elements at the BS  $N_t$  ranges from 100 to 500. The total transmit power of the MIMO-BS is 40 dBm. The transmit power of the APs is 33 dBm. The train has six carriages with

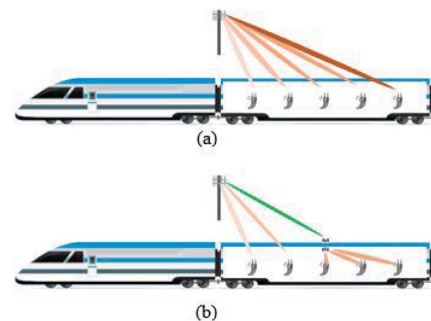
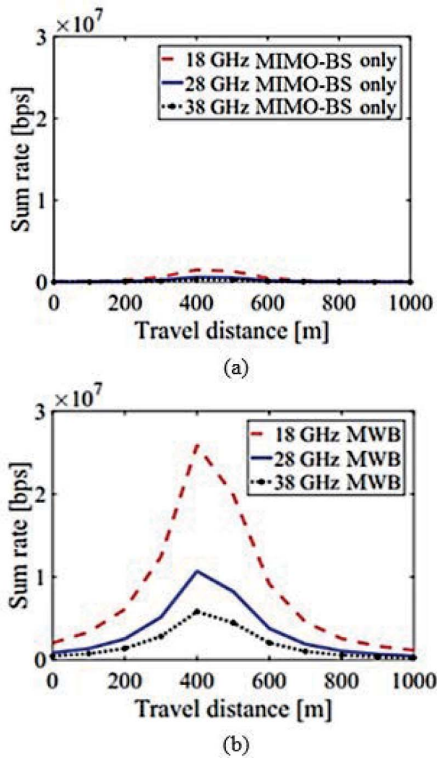


Fig. 23. Considered architectures. (a) MIMO-BS only. (b) MWB.



**Fig. 24.** Influence of train position on sum rate of UEs. (a) The result of the MIMO-BS only architecture; (b) the result of using the MWB architecture.

80 randomly distributed UEs per carriage. The penetration loss of the wagon is 15 dB at all the three frequencies. The length of each carriage is 20 m, and the total length of the train is 120 m. The train moves along the  $x$ -axis from 0 to 1000 m, and the MIMO-BS is located at  $x = 400$  m, which is half the traveling distance. The 2-D horizontal distance  $d_1$  of the MIMO-BS to the track is parallel with the  $y$ -axis, and  $d_1$  ranges from 10 to 100 m. The heights of the MIMO-BS and APs are 20 and 4 m, respectively.

Fig. 24 shows the effect of the train position on the sum rate of UEs ( $N_t = 500$ ,  $d_1 = 10$ ). Fig. 24(a) shows the result of the MIMO-BS only architecture whereas Fig. 24(b) shows the result of using the MWB architecture. In both architectures, the maximum data rate is achieved when the center of the train reaches  $x = 500$  m because the average distance from the APs to the MIMO-BS ( $d_2$ ) is the smallest. The larger the  $d_2$  is, the less the data rate would be. The data rate of the MWB architecture is higher than that of the MIMO-BS-only architecture. When the frequency increases, the path loss increases, and the data rate decreases.

### G. URLLC for Smart Rail Systems

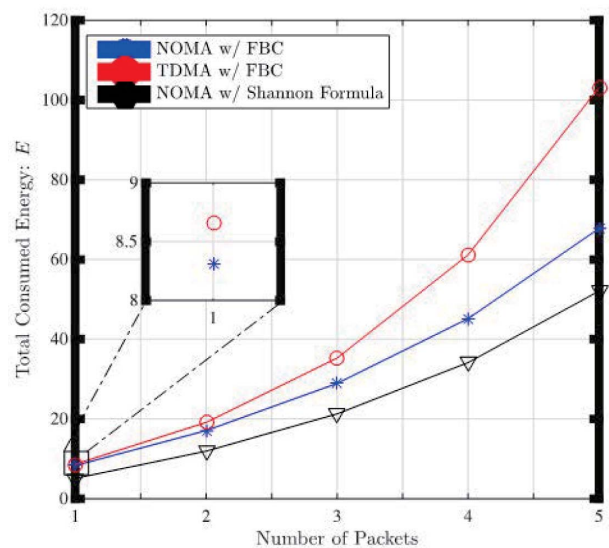
The core performance metric of URLLC is that the packet error ratio is on the order of  $10^{-5}$  under less than 1-ms E2E transmission delay, which is significant for the future development of smart rail communication systems [136].

As one of 5G key technologies, URLLC technology can offer super highly reliable and low E2E latency communication for HSR services, such as safety and control message delivery and surveillance. However, the classical Shannon capacity is accurate only when the codeword blocklength is infinitely long, and thus not applicable to the URLLC systems. Recently, the achievable rate with FBCs over the AWGN channel has been accurately approximated. With received SNR  $\rho$ , blocklength in units of symbols  $m$ , and packet error probability (PEP)  $\epsilon$ , the transmission rate in bits per channel use is given by [137]

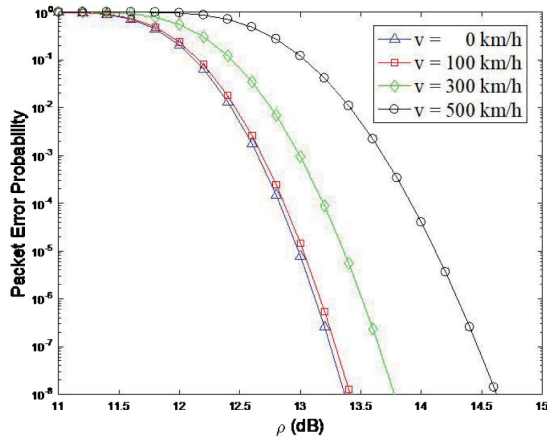
$$\frac{N}{m} = \log_2(1 + \rho) - \sqrt{\frac{1}{m} \left(1 - \frac{1}{(1 + \rho)^2}\right)} \frac{Q^{-1}(\epsilon)}{\ln 2} + \frac{\mathcal{O}(\log m)}{m} \quad (10)$$

where  $N$  denotes the number of bits to be sent. This formula explicitly characterizes the tradeoff between the transmission rate, latency, and reliability. Meanwhile, the theoretical rate of (5) can be practically approached via polar codes with short blocklength. A more recent result given in [138] shows that the rate gap is only 0.025 dB by a simple CRC-polar concatenated scheme with the CRC-aided hybrid decoding method. Hence, the FBC-based formula has been applied to the study of various communication scenarios with strict latency constraints, as in [139].

In Fig. 25, we compare the predicted energy consumptions by invoking the FBC capacity and Shannon capacity of a downlink single-antenna system, in which the energy-efficient resource allocation for a two-user heterogeneous NOMA downlink with an FBC is investigated. The packets may not be aligned thus the successive interference



**Fig. 25.** Energy consumption using FBC and Shannon capacities versus the number of packets, with  $m_1 = m_2 = 640$ ,  $P_{max} = 40$  dBm, and 32 bytes/packet.



**Fig. 26.** PEP  $\varepsilon$  versus received SNR  $\rho$  at different moving speeds, in which the ICI due to the Doppler spread is estimated by its upper bound,  $N = 128$  bytes,  $m = 256$  symbols,  $f_c = 800$  MHz, and  $T_S = 160 \mu\text{s}$ .

cancellation (SIC) cannot be always perfectly performed, which is a key challenge for NOMA transmission with finite blocklength coding. By fixing the requirements of latency ( $m_1 = m_2 = 640$ ) and PEP  $\varepsilon = 10^{-6}$ , it shows that the Shannon capacity formula universally underestimates the energy of both the TDMA and NOMA schemes. We can also observe that the performance gains of NOMA, as compared with those of TDMA, increases as the number of transmitted packets increases due to the higher SE from the nonorthogonal superposition coding. This indicates that the performance caused by FBC should be carefully considered in URLLC system designs. Note that in the simulation, we assume QAM modulation with channel code 1/5, the system bandwidth is 1 MHz and the noise power density is set to be  $\sigma_1^2 = \sigma_2^2 = -110$  dBm. The interested readers are referred to the details to [139]. Note that the PEP performance can be further improved by adaptive coding.

We can also observe that the performance gains of NOMA, as compared with those of TDMA, increase as the number of transmitted packets increases due to the higher SE from the nonorthogonal superposition coding. This indicates that the performance caused by FBC should be carefully considered in URLLC system designs.

Meanwhile, the FBC capacity degrades significantly under the imperfect CSI. Specifically, in high-speed scenarios, the CSI estimation is challenging even in a typical case where trains run over viaducts; thus, the channel is dominated by the LoS with a large coherence bandwidth. There are three main reasons for this: 1) fast-varying CSI due to high mobility; 2) the ICI of an OFDM system due to the Doppler spread; and 3) limited number of pilot symbols and feedback delay.

In Fig. 26, we compare the PEP performance versus received SNR under different movement speeds. To simplify, we assume that the received SNR  $\rho$  is known whereas

the ICI power is approximately characterized by its upper bound

$$P_{\text{ICI}} \leq \frac{1}{12}(2\pi f_d T_S)^2 \quad (11)$$

where  $f_d = v f_c / C$  denotes the maximal Doppler frequency,  $f_c = 800$  MHz,  $v$  is the movement speed,  $C = 3.0 \times 10^8$  m/s, and  $T_S = 160 \mu\text{s}$ . In the case with  $N = 128$  bytes and  $m = 256$  symbols, we find that the PEP is very sensitive to the received SNR; specifically, it increases as the movement speed increases. If the speed increases from 300 to 500 km/h, then it would require 1-dB SNR gain to achieve the same PEP of  $10^{-6}$ . This implies that both the channel estimation and the resource allocation are two key issues of the URLLC system design for high-speed applications.

## H. Power Adjustment Enhanced Handover for Smart Railways

Many existing works showed that increasing transmit power is capable of reducing the handover failure probability, but it is power hungry and goes against the indicator of green communications in 5G. Thus, adjusting the power in the time domain may be an efficient way, since it is able to reduce the “uncertainty” of the received signal strength in the handover procedure of HSR communications and does not require more transmit energy during the handover procedure [140], [141]. The basic idea is described as follows.

Consider an HSRC system, where a train is moving through an overlap region of two neighboring cells. The moving users within the train desire to connect with the ground BS via the TRS and some APs installed inside the carriages over a two-hop, that is, user-TRS-BS, communication architecture. When moving users intend to access the Internet via the BS, in the uplink, APs first collect the data from the users and then forward the aggregated data and requests to the BS with the help of TRS. In the downlink, the BS first transmits information to the TRS and then the TRS helps to forward the information to users via the APs. With such a two-hop architecture, as mentioned previously, the heavy overhead of the group handover is mitigated.

To perform the power adjustment, the overlap region is divided into two subregions, as shown in Fig. 27. The first one is from point A to point B, and the second one is from point B to point C, where A and C denote the starting and the ending points of the overlap region, respectively. B is the middle point. In the first subregion, the transmit power of the serving cell is increased and that of the target cell is decreased. In the second subregion, vice versa. To keep the energy consumption of each BS unchanged within the overlap region, the power is adjusted such that the increased amount of energy is equal to the decreased amount. Such a location-based time-domain power adjustment can be realized by estimating

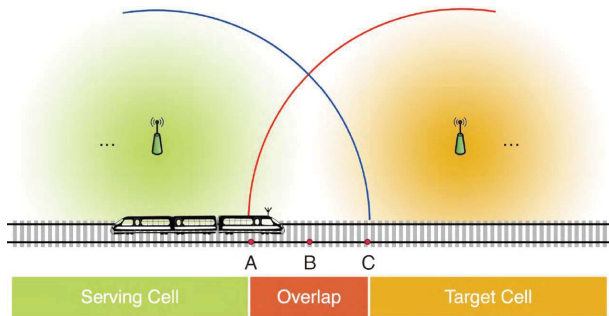


Fig. 27. Illustration of the subregions associated with handover.

Table 4 Simulation Parameters

Parameter	Value
Current transmit power of BS	46dBm
Maximum transmit power of BS	50dBm
Adjusted power	4dB
Shadow fading deviation	4dB
Path loss model	$31.5 + 35\log_{10}(d)$
Distance between two cells	3000m
Number of RAUs in one cell	4
Noise density	-145dBm/Hz
Signal threshold	-30dB
Predefined margin	2dB
Distance from BS to railways	100m
Distance from RAU to railways	60m
Speed of train	300m/s

the channel information based on the position of the train.

To show the handover performance of the power adjustment enhanced method, we simulate and also compare it with other benchmarks with the following parameters as shown in Table 4.

In Fig. 28, both the analytical and the simulation results are provided, where the curves marked with “with PA” and “w/o PA benchmark x” represent the performance of the handover methods with and without power adjustment, respectively. Particularly, benchmark 1 is the handover scheme in traditional cells, which is with no RAU and PA. Benchmark 2 is the blanket transmission-based handover

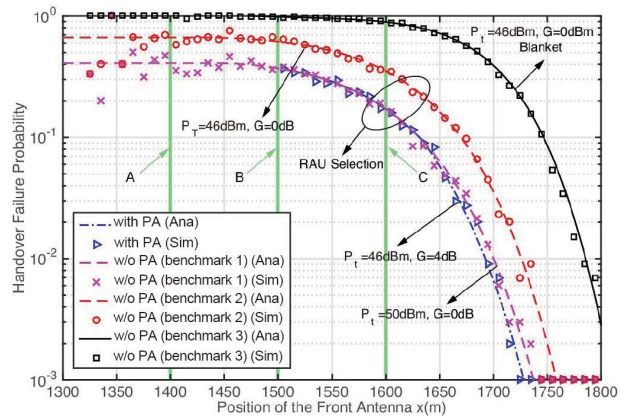


Fig. 28. Handover failure probability.

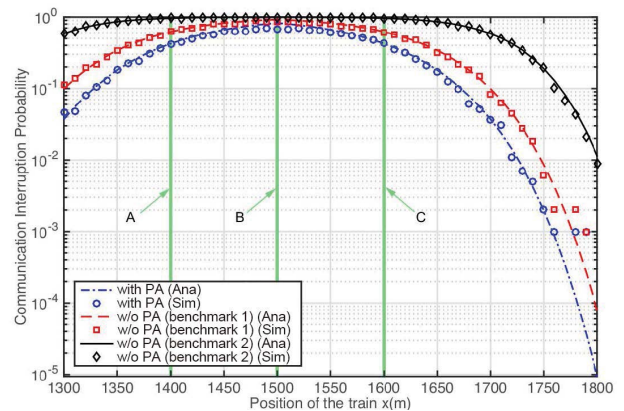


Fig. 29. Communication interruption probability.

scheme in DAS cells without PA. Benchmark 3 is the RAU selection transmission-based handover scheme in DAS cells without PA. It is shown that by including power adjustment into handover, a much lower handover failure probability is achieved within the subregion x from B to C [140]. With power adjustment, the communication interruption probability is greatly decreased in the handover procedure [140], which is shown in Fig. 29.

## VII. FUTURE WORKS

In this article, we have provided an overview of 5G-R emphasizing network architectures, channel models, and key technologies for future smart rails. However, there are still many challenges before we can successfully deploy 5G networks and services in the future railway scenarios. Such railway scenarios include D2D, integrated ground-air-space network, IoT, security, and AI, which are considered as key emerging technologies in 5G networks. In this part, we present some key new technology components and solutions to address the key challenges and requirements for a future smart railways.



### A. D2D for Smart Railways

When the wireless link between the BS and the train fails, operators could consider direct D2D communication as an emergency communication tool that is a well-integrated part of the 5G overall wireless-access solution [142]. Using direct D2D communications could be considered as a means of extending coverage beyond the reach of the conventional infrastructure (device-based relaying). D2D can also provide infrastructure fault-alerting and monitoring mechanism inside train carriages.

Applying D2D to railway scenarios is challenging for several reasons. First, D2D communications should also be possible in scenarios where network coverage is unavailable. As such, a D2D link needs to establish connectivity without- network control/assistance. Second, the D2D link is unstable because trains are highly mobile.

Cooperative D2D have high-speed communications that provides operators with the means for “joint” transmission and/or reception between multiple devices. Note that this can be seen as a kind of coordinated transmission and reception; however, this would be on the device side and not on the network side.

The key point is to see D2D communications as a well-integrated part of the overall wireless access solution. Accordingly, D2D communications should be considered right from the start of the 5G definition, instead of taking it as a later-introduced “add-on” option.

### B. Integrated Ground-Air-Space Network for Smart Railways

Using satellites, airships, unmanned aerial vehicles, and other air platforms in 5G-R scenarios has many advantages [143]. For example, it can provide efficient E2E transmission services for data, such as smart rail network-detection data, early warning information, and remote decision-making. However, achieving integrative information transmission and processing under air, space, train, and ground settings is challenging for 5G-R.

First, a high data-rate wireless connectivity should be built to support the tremendous amount of data that the integrated ground-air-space network will produce. Second, power-limited airships and unmanned aerial vehicles cannot be used under strong rain or wind conditions. To solve these problems, designers should provide an efficient data distribution and resource-sharing scheme under a unified control.

### C. IoT for Smart Railways

In addition to real-time query and tracking the whole trajectory of the train and goods location, the IoT for railways can be developed in 5G-R networks to integrate sensing information of rail infrastructures, including bridges, viaducts, tunnels, leaky feeders, rail gaps, frozen soil, and slope protection [144]. This could be done by installing various sensing measures such as infrared sensors, sound sensors, and temperature sensors.

In recent years, IoT technology has gradually attracted the attention of railway departments, industries, and research institutes all over the world. Developing railway IoT and building railway safety information-guaranteed systems based on IoT are important directions for the in-depth integration of railway information and industrialization.

However, researchers should further explore the mobile communication problems in IoT, and then conduct research on multiple access from a massive number of devices in the presence of high mobility. Likewise, the mechanisms for evaluating and optimizing railway wireless network resource management mechanism should also be further studied. Designers also need to develop and construct a closed-loop management of safety monitoring and control system for railway equipment and facilities, collect real-time information of mobile equipment and fixed facilities, analyze the disposition, and ensure safe operation. Also, monitoring systems for natural disasters (e.g., wind, rain, snow, and earthquakes) and foreign matter intrusion need to be established to monitor the safety of railway operations in real time.

### D. Security for Smart Railways

The operator of the 5G-R system should be able to securely collect information that can enhance user and service experiences via data analytics. Security has been one of the fundamental capabilities that operators should provide to their customers, especially to train operation control systems [145]. 5G-R will support a wide range of applications, from human-based to machine-based communications. Thus, the corresponding technology should be able to deal with huge amounts of sensitive data that need to be protected against unauthorized access, use, disruption, modification, inspection, attack, etc. The importance of providing a comprehensive set of features that guarantee high-level security is a core requirement for 5G-R systems. Therefore, 5G-R should be designed to provide more options beyond node-to-node and E2E security that is available in today’s mobile systems. Accordingly, such a design will protect the users’ data and prevent or mitigate any possible cyber security attack.

On the other hand, operators should protect railway communications and signaling against electromagnetic attacks, which can identify scenarios and devices that are vulnerable to attacks. Likewise, researchers and developers should conduct a risk analysis of attack scenarios, propose responses, develop solution to detect electromagnetic attacks, and design the right architecture to be resilient to such kind of attack.

### E. AI for Smart Railways

Recently, AI has become popular in the computer field due to its great success in computer vision, natural language processing, automatic speech recognition, and wireless communications [146]. Future smart rail considers

AI as a significant direction in 5G-R networks; AI would enable networks to process a large amount of data, dynamically recognize and adapt to complex scenarios, and satisfy the requirements for high-speed and real-time signal processing capabilities.

Note that most AI methods existing today are data driven. They use the standard neural network structure as a black box and then train it through large amounts of data. By integrating AI into the physical and upper layers, smart rail communication systems would be able to automatically adapt to the properties of signals and propagation scenarios, and thereby obtain efficient deployment. However, the current design of AI schemes for communication networks is simple. To improve the network performance, we should resort to a specialized AI architecture for 5G-R communications that considers the special characteristics of railway scenarios and strict indicators. Above all, with the development of advanced communication technologies, future railways will become smarter and smarter.

## F. Mobile Edge Computing (MEC) for Smart Railways

Under HSR scenarios, the high-speed movement of the train faces the dramatically time-variance of the channel and the rapid reduction of the available bandwidth. The new applications for smart rail, such as automatic train driving, real time management for high-speed rail networks, and real time HD video surveillance, need stronger computing capabilities and lower processing latencies.

MEC is considered as an emerging key technique in 5G communication which extends the ability of cloud computing to the edge of the network and reduces the bandwidth requirement between the train and the BS [147]. MEC for smart rail consists of onboard MEC and trackside MEC. Onboard MEC stores and processes data that are incoming with a high rate (HD video surveillance data, train sensor data) generated during train running. The processed results, which require a small data rate, are fed back and shared with trackside MEC. Trackside MEC analyzes and processes the data of wayside sensors, receives data from onboard MEC, and shares the processed data with the cloud computing node. With the deployment of MEC, the process of smart rail data will become real-time and efficient.

## G. Sixth-Generation (6G) for Smart Railways

As 5G is deployed around the world, 6G becomes a hot topic and has attracted increasing attention from researchers in both academia and industry [148]. The development of 6G will pave the way for diverse QoS requirements, real-time tactile interactions, customized and open services, integration of communication, broadcasting [149]–[151], computing, sensing, control, security, and AI functionalities. The main innovation in 6G network architecture is also required with respect to current 5G

network design, such as space-aerial-terrestrial-sea integrated network, full-spectrum and full-dimensional coverage, intelligent self-sensing, learning, optimization, and evolution.

On the other hand, the potential applications for future smart railway networks include autonomous train driving, cooperative train networks, Internet of trains, UHD (4K/8K) train video, train *ad hoc* networks, and ultraaccurate (cm level) train localization. In order to fulfill the requirements of 6G smart railway applications, significant technological innovations are expected.

- 1) *Cell-Free network architecture*: The high mobility of trains induces frequent handover and huge overhead in conventional cellular networks. In order to guarantee a seamless and high-QoS coverage, it is promising to employ cell-free networks. Moreover, different heterogeneous sub-6G, mmWave, terahertz, and visible light communications enable cell-free networks to provide unprecedented performance. The application of AI techniques in 6G smart railway networks is promising in the context of self-learning, optimization, and evolution.
- 2) *Novel URLLC techniques*: To guarantee the autonomous train driving over 1000 km/h in 6G (e.g., emerging Hyperloop systems), it is important to further investigate the novel URLLC frame structure under FBC to achieve the tradeoff between ultrareliability and ultralow latency.
- 3) *Digital twin networks*: The monitoring, prediction and decision-making for the behavior of train drivers (e.g., nervous, drunk, sleepy, and excited) is significant for the security of HSTs. A digital twin network based on 6G can fundamentally evolve the digital profile of the historical and current behavior of the train driver that helps optimize their performance.

Finally, we believe the networking, broadcasting, communications, interaction, and security of future smart railways will be vastly improved by applying the aforementioned techniques, which deserve our effort in the next decade.

## VIII. CONCLUSION

Wireless communication is playing an important role in the future smart railway domain. This article explored a potential solution by leveraging emerging 5G technologies to provide a plethora of services in HSRs, both control and data services. More specifically, we first briefly described the current trend of wireless communications for smart railways. Services and requirements of future smart railways were presented in the following. Moreover, we introduced state of the art, drawbacks, and challenges of existing wireless technologies in supporting the envisioned services, respectively. We proposed a new 5G-R network architecture for smart railways. Channel models, physical layer design challenges and a thorough review







## ABOUT THE AUTHORS

**Bo Ai** (Senior Member, IEEE) graduated from Tsinghua University, Beijing, China, with the honor of Excellent Postdoctoral Research Fellow in 2007. He received the master's and Ph.D. degrees from Xidian University, Xi'an, China, in 2002 and 2004, respectively.



He was a Visiting Professor with the Electrical Engineering (EE) Department, Stanford University, Stanford, CA, USA, in 2015. He is currently working as a Full Professor and a Ph.D. Advisor with Beijing Jiaotong University, Beijing, where he is also the Deputy Director of the State Key Laboratory of Rail Traffic Control and Safety and the International Joint Research Center. He is one of the main responsible people for Beijing (urban rail operation control system) International Science and Technology Cooperation Base, and the Backbone Member of the Innovative Engineering-Based jointly granted by the Chinese Ministry of Education and the State Administration of Foreign Experts Affairs. He has authored or coauthored eight books and published over 300 academic research articles in his research area. He holds 26 invention patents. He has been the research team leader for 26 national projects and has won some important scientific research prizes. Five of his articles have been the ESI highly cited articles. He has been notified by the Council of Canadian Academies (CCA) that, based on Scopus database, he has been listed as one of the Top 1% authors in his field all over the world. He has also been feature interviewed by IET Electronics Letters. His research interests include the research and applications of channel measurement and channel modeling and dedicated mobile communications for rail traffic systems.

Dr. Ai is a Fellow of the Institution of Engineering and Technology (IET). He is also an Editorial Committee Member of *Wireless Personal Communications* journal. He has received many awards, such as the Distinguished Youth Foundation and Excellent Youth Foundation Award by the National Natural Science Foundation of China, the Qushi Outstanding Youth Award by the Hong Kong Qushi Foundation, the New Century Talents by the Chinese Ministry of Education, the Zhan Tianyou Railway Science and Technology Award by the Chinese Ministry of Railways, and the Science and Technology New Star Award by the Beijing Municipal Science and Technology Commission. He is also a Distinguished Lecturer of the IEEE Vehicular Technology Society, a Vice Chair of the IEEE VTS Beijing Chapter, and the Chair of the IEEE BTS Xi'an Chapter. He was the co-chair or the session/track chair for many international conferences. He is also an Associate Editor of the IEEE ANTENNAS AND WIRELESS PROPAGATION LETTERS and the IEEE TRANSACTIONS ON CONSUMER ELECTRONICS. He is the Lead Guest Editor for Special Issues on the IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY, the IEEE ANTENNAS AND PROPAGATIONS LETTERS, and the *International Journal on Antennas and Propagations*.

**Andreas F. Molisch** (Fellow, IEEE) received the Dipl.Ing., Ph.D., and Habilitation degrees from the Vienna University of Technology, Vienna, Austria, in 1990, 1994, and 1999, respectively.



He spent the next ten years in industry at FTW, Austria, AT&T (Bell) Laboratories, USA, and Mitsubishi Electric Research Laboratories, Cambridge, MA, USA, where he rose to the Chief Wireless Standards Architect. In 2009, he joined the University of Southern California (USC), Los Angeles, CA, USA, as a Professor, where he founded the Wireless Devices and Systems (WiDeS) Group. In 2017, he was appointed as the Solomon Golomb–Andrew and Erna Viterbi Chair. Overall, he has published four books (among them the textbook *Wireless Communications* currently in

its second edition), 20 book chapters, 250 journal articles, and 350 conference papers. He is also an inventor of 60 granted (and more than 20 pending) patents and a coauthor of some 70 standard contributions. His research interests revolve around wireless propagation channels, wireless systems design, and their interaction. Recently, his main research interests include wireless channel measurement and modeling for 5G and beyond 5G systems, wireless video distribution, hybrid beamforming, ultrawideband/time-of-arrival-based localization, caching at the wireless edge, and novel modulation/multiple access methods.

Dr. Molisch is a Fellow of the National Academy of Inventors, the American Association for the Advancement of Science, and the Institution of Engineering and Technology (IET). He is also a member of the Austrian Academy of Sciences. He has been an editor of a number of journals and special issues, the general chair, the technical program committee chair, or the symposium chair of multiple international conferences, and the chairman of various international standardization groups. He is also an IEEE Distinguished Lecturer. He has received numerous awards, among them the IET Achievement Medal, the Technical Achievement Awards of the IEEE Vehicular Technology Society (Evans Avant-Garde Award), the IEEE Communications Society (Edwin Howard Armstrong Award), the Technical Field Award of the IEEE for Communications, and the Eric Sumner Award.

**Markus Rupp** (Fellow, IEEE) received the Dr.Ing. degree from the Technische Universität Darmstadt, Darmstadt, Germany, in 1993.



He served as the Dean from 2005 to 2007 and 2016 to 2019 and the Head of the Institute of Telecommunications, Technische Universität Wien, Vienna, Austria, from 2014 to 2015, where he has been a Full Professor since 2001. He has authored or coauthored more than 650 scientific articles and patents on adaptive filtering, wireless communications, rapid prototyping, and automatic design methods.

Dr. Rupp was with the Board of Directors (BoD) of EURASIP from 2004 to 2011, including two years serving as the President of the society.

**Zhang-Dui Zhong** (Senior Member, IEEE) is currently a Professor and an Advisor of Ph.D. candidates with Beijing Jiaotong University, Beijing, China. He is also the Director with the School of Computer and Information Technology and a Chief Scientist with the State Key Laboratory of Rail Traffic Control and Safety, Beijing Jiaotong University. He is also the Director of the Innovative Research Team, Ministry of Education, and a Chief Scientist of the Ministry of Railways, Beijing. He has authored or coauthored seven books, five invention patents, and more than 200 scientific research articles in his research area. His research interests include wireless communications for railways, control theory and techniques for railways, and global system for mobile communications for railway (GSM-R) systems. His research has been widely used in the railway engineering, such as Qinghai-Xizang Railway, Datong-Qinhuangdao Heavy Haul railway, and many high-speed railway lines of China.



Dr. Zhong is an Executive Council Member of the Radio Association of China and the Deputy Director of the Radio Association of Beijing. He received the Mao Yisheng Scientific Award of China, the Zhan Tianyou Railway Honorary Award of China, and the Top Ten Science/Technology Achievements Award of Chinese Universities.