

Toward Federated Cognitive Digital Twins over the Edge-to-Cloud Continuum

Alessandra Somma  

University of Naples Federico II, Italy

Alessio Bucaioni  

Mälardalen University, Sweden

Abstract

Digital Twins (DTs) are increasingly adopted to monitor, analyze, and optimize Cyber-Physical Systems (CPSs) by enabling continuous interaction between physical assets and their digital representations. Despite their potential, current DT architectures face significant limitations when applied to distributed environments such as smart cities. In particular, most existing solutions rely on centralized and monolithic designs, which introduce latency, scalability, and resilience issues. Moreover, they provide limited support for semantic integration and high-level reasoning, which hinders the effectiveness of DTs-based decision-making.

Recent research has explored Federated Digital Twins (FDTs) to decompose complex systems into multiple interacting twins, improving scalability and modularity. However, existing FDT approaches often lack a clear architectural framework and still concentrate intelligence in cloud-based components, thus limiting local autonomy. In parallel, Cognitive Digital Twins (CDTs) have been proposed to enhance DTs with semantic reasoning and explainability, leveraging Artificial Intelligence and Large Language Models (LLMs). While promising, these approaches are typically centralized and difficult to integrate within distributed architectures.

In this paper, we propose a unified approach that combines federation and cognition within a single architecture, referred to as the Federated Cognitive Digital Twin (FCDT). The proposed architecture distributes intelligence across the edge-to-cloud continuum by introducing local twins deployed close to physical systems and global twins operating at the cloud level. Local twins provide real-time monitoring, analysis, and first-level decision-making, enhanced by lightweight cognitive capabilities; global twins perform system-level reasoning, simulation, and coordination, and leverage more computationally intensive models for explanation and decision support.

By combining distributed autonomy at the edge with global cognitive simulation and reasoning in the cloud, the FCDT architecture addresses both structural and semantic limitations of current DT solutions. It enables scalable and responsive DT systems, improves decision-making capabilities, and simplifies the engineering of complex CPSs in distributed environments.

2012 ACM Subject Classification Software and its engineering → Software design engineering; Computer systems organization → Cloud computing

Keywords and phrases Digital Twins, Cloud Continuum, Edge Computing, Federation, LLMs

Digital Object Identifier 10.4230/OASICS.CVIT.2016.23

1 Introduction

Digital Twins (DTs) are virtual representations of Cyber-Physical Systems (CPSs), characterized by continuous bidirectional communication between the physical and digital domains [7]. This interaction enables the transfer of real-world data from the physical system to its digital counterpart, while also allowing the DT to generate insights and provide feedback to the physical world [4, 25]. By combining data, models, and analytics, DTs support advanced monitoring, enhance situational awareness, and enable data-driven decision-making across domains such as manufacturing, transportation, and smart cities [2, 23].

Engineering and architecting DT is a complex and demanding activity, and for this reason



© Alessandra Somma and Alessio Bucaioni;
licensed under Creative Commons License CC-BY 4.0
42nd Conference on Very Important Topics (CVIT 2016).

Editors: John Q. Open and Joan R. Access; Article No. 23; pp. 23:1–23:10

OpenAccess Series in Informatics



OASICS Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

46 it has received growing attention in recent years. This interest has led to the proposal of
47 many DT architectures [7,15], together with important standardization efforts [6,9]. However,
48 when DTs are applied to large-scale and inherently distributed environments such as smart
49 cities, new challenges emerge [23]. In these contexts, heterogeneous devices continuously
50 generate data, multiple subsystems interact with one another, and decisions often need to be
51 taken in a timely and coordinated way. As a result, current DT approaches reveal two major
52 weaknesses: they are still largely based on centralized designs [5,7,24], and they provide
53 limited support for semantic integration across heterogeneous subsystems and information
54 sources [8,24].

55 Most existing DT architectures address this complexity through centralized solutions. In
56 centralized designs, the DT is usually implemented as a monolithic system, often deployed in
57 the cloud, and designed as a single unit representing the whole real system. At the same time,
58 data generated by distributed physical assets are transmitted to centralized platforms, where
59 are then stored and processed to feed simulation models, analytics pipelines, and AI-based
60 inference engines. This design has clear advantages: it simplifies DT management, supports
61 large-scale data processing, and enables advanced monitoring and simulation capabilities.
62 However, these benefits come at a cost. The concentration of both the DT logic and the data
63 processing pipeline in centralized infrastructures introduces significant challenges in terms of
64 latency, bandwidth consumption, and resilience [20,21]. Centralized solutions remain also
65 sensitive to network disruptions and make local autonomy difficult to achieve [17].

66 For these reasons, DT research is progressively moving beyond monolithic and fully
67 centralized designs toward more distributed solutions. In this direction, the concept of
68 **Federated Digital Twins** (FDTs) has emerged as a promising way to decompose the
69 complexity of large CPSs [26,27]. The main idea is to model the overall system as a set
70 of autonomous yet interoperable twin instances, each associated with a specific subsystem,
71 while allowing them to cooperate in order to provide a coherent system-level view [16,27]. In
72 this way, federation addresses an important structural limitation of monolithic DTs, since
73 the complexity of the global system is no longer handled by a single twin, but distributed
74 across multiple coordinated twins.

75 However, the FDT paradigm is still in an early stage, and many existing proposals remain
76 at a rather high level of abstraction. In particular, they often lack a clear architectural
77 framework to support coordinated intelligence across distributed twins [14,23,28]. Moreover,
78 even when federation is adopted, intelligence is frequently still concentrated in cloud-based
79 components. This helps in managing the overall complexity of the DT, but it does not
80 fully solve the problems that motivated the shift away from centralization in the first place,
81 especially latency, resilience, and local autonomy. The result is that currently federation of
82 DTs improves their structural decomposition, but often without yet providing a fully effective
83 distribution of their intelligence [14,28].

84 Hence, the shift toward federation helps in addressing the structural complexity of
85 distributed CPSs [13,27], but it does not fully solve their limited semantic capabilities. As a
86 consequence, DTs are often able to detect or predict events, but they struggle to explain
87 them, relate them across subsystems, or support reasoning about their implications. For
88 example, identifying a delay in a transportation network does not automatically provide an
89 understanding of its causes, how it propagates through the system, or which actions should
90 be taken to mitigate its impact. This gap between data processing and meaningful reasoning
91 limits the effectiveness of DTs as decision-support systems [12].

92 To address this limitation, recent research has introduced the concept of **Cognitive**
93 **Digital Twins** (CDTs) [12,31], which extend traditional DTs with capabilities for semantic

94 reasoning, knowledge integration, and explainability. The idea is to move from DTs that only
95 process data to DTs that can also interpret it, relate it to domain knowledge, and support
96 human understanding [1, 30]. In this context, advances in Artificial Intelligence (AI), and
97 in particular in Generative AI (GenAI), play a key role. Large Language Models (LLMs)
98 have shown strong capabilities in integrating heterogeneous information, reasoning across
99 domains, and interacting with users through natural language [12, 29, 32].

100 However, introducing this type of intelligence also raises new challenges. LLMs typically
101 require significant computational resources and are commonly deployed in centralized environ-
102 ments. This creates a tension with the need for distributed DT architectures, especially in
103 scenarios where low latency and local autonomy are required [22]. If not carefully designed,
104 the integration of such models risks reinforcing centralization, which is precisely one of the
105 limitations that federated DTs aim to overcome.

106 Recent developments suggest that this tension can be mitigated by distributing intelligence
107 across the edge-to-cloud continuum [3, 10, 11, 19]. In this setting, lightweight AI models,
108 including distilled or specialized versions of language models, can be deployed at the edge
109 to support local tasks such as filtering, summarization, and contextual interpretation [18].
110 At the same time, more computationally intensive models can be deployed in the cloud to
111 provide global reasoning, explanation, and decision support [3, 10].

112 In this paper, we propose to address these challenges by proposing a unified architecture
113 combining federation and cognition, resulting in the **Federated Cognitive Digital Twin**
114 (FCDT) architecture. In this architecture, local twins provide fast and autonomous decision-
115 making close to the physical system, while global twins support system-level reasoning and
116 coordination. By combining these two levels, the proposed approach maintains responsiveness
117 while enabling advanced analysis and explanation, and at the same time simplifies the
118 engineering of complex DT systems and mitigates latency issues in distributed environments.

119 **2 Related Work**

120 Existing research on engineering and architecting DTs has evolved along multiple directions,
121 addressing different aspects of DTs design, distribution, and intelligence. In this section,
122 we focus on two lines of research, which are relevant to our work, i.e., Federated DTs and
123 Cognitive DTs.

124 Early DT architectures are predominantly based on centralized designs, where data from
125 distributed assets are collected and processed in cloud environments [5, 7, 15, 24]. FDTs
126 have been introduced to address the structural limitations of monolithic DT architectures
127 in large-scale systems. Instead of representing the entire real system through a single DT
128 instance, FDT approaches decompose it into multiple interacting twins, each associated with
129 a specific subsystem.

130 Vergara *et al.* [26, 27] define FDTs as a distributed set of twin instances that collaborate to
131 support system-level decision-making. Their work emphasizes the role of federation in enabling
132 collaborative simulation and coordination, where different twins exchange information to
133 improve global understanding. Papacharalampopoulos *et al.* [16] analyze federation from
134 a modeling perspective, highlighting limitations in knowledge transfer and interoperability
135 between twin instances. Their work shows that even when multiple twins are connected,
136 inconsistencies and lack of shared semantics can hinder effective collaboration. This reinforces
137 the need for stronger integration mechanisms beyond simple federation.

138 Model-driven approaches, such as the one proposed by Michael *et al.* [14], investigate
139 how FDTs can be engineered using abstraction and modeling techniques. These approaches

140 contribute to structuring the development process, but they still operate at a high level and
141 do not explicitly address runtime coordination or distributed decision-making.

142 More recent efforts, such as [28], explore learning-based coordination across federated
143 twins, while large-scale initiatives like interTwin [13] combine federation with distributed
144 computing infrastructures. These works show the potential of FDTs in complex environments,
145 but they also reveal a common limitation: intelligence is often still centralized, typically in
146 cloud-based components that aggregate and process data from all twins.

147 As a result, current FDT approaches successfully address structural decomposition and
148 scalability, but only partially address distributed intelligence. In many cases, federation is
149 used to organize the system, but decision-making and reasoning remain centralized, which
150 limits responsiveness and local autonomy.

151 In parallel, CDTs have been proposed to address the limited reasoning capabilities of
152 traditional DTs. CDTs extend DTs with semantic reasoning, knowledge integration, and
153 explainability, moving beyond purely data-driven analytics [12,31]. Zheng *et al.* [31] introduce
154 the concept of cognitive DTs as systems that integrate data, models, and knowledge to support
155 decision-making. Their work highlights the importance of combining physical models with
156 higher-level reasoning mechanisms, but does not provide a concrete architectural realization.

157 Early CDTs works such as [1] and [30] propose hybrid DT architectures that combine data-
158 driven models with symbolic reasoning and self-awareness mechanisms. These approaches
159 represent an important step toward cognitive DTs, but their scope is often limited to specific
160 domains and they do not address large-scale distributed systems.

161 Liu *et al.* [12] present a comprehensive survey on CDTs and discuss how AI techniques can
162 enhance DT capabilities. In particular, they identify the role of knowledge graphs, reasoning
163 engines, and more recently LLMs in enabling semantic understanding and interaction. Indeed,
164 with the rise of Generative AI, several works have explored the integration of LLMs into
165 DT systems. For instance, Xia *et al.* [29] propose an architecture where LLMs are used
166 to interact with DTs and support reasoning over system data. Similarly, Zhu *et al.* [32]
167 discuss how semantic communication and LLM-based reasoning can enhance DT systems.
168 These works demonstrate the potential of LLMs in enabling explainability, interaction, and
169 cross-domain reasoning.

170 Recent research has started to explore how AI models, including LLMs, can be distributed
171 across the edge-to-cloud continuum. Surveys such as [10] highlight the potential of deploying
172 lightweight models at the edge to support local processing, while keeping more complex
173 models in the cloud.

174 Frameworks such as Edge-LLM [3] and hybrid inference approaches [11] propose splitting
175 LLM execution across edge and cloud resources. These approaches enable tasks such as
176 filtering, summarization, and contextual interpretation to be performed locally, reducing
177 communication overhead and improving responsiveness. Similarly, [18] explores optimization
178 strategies for edge-cloud LLM deployment.

179 These approaches suggest that LLM-based capabilities can be adapted to distributed
180 environments through techniques such as model distillation, specialization, and collaborative
181 inference. In this setting, edge components handle lightweight tasks, while cloud components
182 perform more complex reasoning and decision-making [18]. Nevertheless, even though these
183 works provide important insights into distributed AI, they are not specifically designed for
184 DT systems. In particular, they do not define how distributed AI components should be
185 integrated within a DT architecture, nor how they should interact with local and global twin
186 instances.

187 The analysis above highlights a clear gap in the current DTs research. Federated

188 approaches address the decomposition of complex systems into multiple twins, but often lack
189 mechanisms for distributing intelligence and coordinating decisions across them. Cognitive
190 approaches enhance reasoning capabilities, but are typically centralized and do not account
191 for the constraints of distributed CPSs. Research on edge AI and LLM distribution provides
192 enabling technologies, but does not offer a DT-specific architectural framework. As a result,
193 existing solutions tend to address either distribution or cognition, but rarely both in a unified
194 way. This gap motivates the need for a new approach that combines federation and cognition
195 within a single architecture.

196 **3 Federated Cognitive Digital Twin Architecture**

197 The proposed FCDT architecture models the Digital Twin as a federated and cognitive
198 system distributed across the edge-to-cloud continuum. The architecture is based on the
199 distinction between *local twins* and *global twins*, which operate at different levels and provide
200 complementary capabilities. Local twins are associated with specific physical subsystems and
201 are deployed close to them at the edge, while global twins operate in the cloud and reason
202 over the information produced by multiple local twins.

203 A key principle of the architecture is that part of the DT capabilities is brought closer to
204 the physical system. This is not only a deployment choice, but a functional one. Local twins
205 are not limited to collecting and forwarding data; they provide monitoring, analysis, and a
206 first level of decision-making directly at the subsystem level. This allows them to maintain
207 tighter synchronization with real-world conditions and to react quickly to local changes.
208 The cloud remains essential, but its role shifts from being the central point of intelligence
209 to the level where information are aggregated, correlated, and used for broader reasoning,
210 simulation, and coordination.

211 Each local twin maintains a continuously updated representation of its subsystem by
212 ingesting data from sensors, devices, and local services. This state forms the basis for real-
213 time analysis, including anomaly detection, short-horizon prediction, and state estimation.
214 Since these activities are performed close to the data source, they can be executed with low
215 latency and without relying on continuous cloud communication, which is crucial in dynamic
216 and time-critical environments.

217 Based on this analysis, local twins support a first level of decision-making. By combining
218 current state and analytical results, they can take immediate actions within the scope of their
219 subsystem. Although these decisions are inherently local, they provide fast and autonomous
220 responses that would be difficult to achieve in centralized systems. In this way, local twins
221 contribute directly to system responsiveness and resilience.

222 In addition to these operational capabilities, local twins incorporate lightweight cognitive
223 components, including edge AI and specialized or distilled edge language models. These
224 components support semantic processing tasks such as event extraction, summarization, and
225 contextual interpretation. This enhances the quality of local decisions by enriching numerical
226 data with contextual meaning, which can be useful for human operators or local control logic.
227 At the same time, it enables local twins to transform raw data into structured and compact
228 representations before sharing it. This process, referred to as semantic distillation, reduces
229 communication overhead while preserving relevant information. It can also support privacy-
230 aware data sharing by enabling the transmission of abstracted or anonymized information.

231 As a result, local twins play a dual role. They act as operational components, responsible
232 for real-time monitoring and local decisions, and as cognitive filters, responsible for trans-
233 forming raw observations into structured knowledge. These outputs are stored in a shared

234 data space and constitute the input for higher-level reasoning.

235 Global twins operate at the cloud level and provide a system-wide perspective. Their role
236 is to reason over the structured information produced by multiple local twins and other global
237 twins, and derive knowledge that cannot be obtained from a single subsystem. Unlike local
238 twins, they do not directly interact with the physical environment, but consume data from
239 the shared data space in the form of events, predictions, and contextualized observations.
240 This abstraction allows them to focus on higher-level and more computationally demanding
241 tasks.

242 Using these aggregated information, global twins perform system-level analysis, including
243 the identification of correlations across subsystems, detection of global patterns, and long-
244 term prediction. They are also responsible for simulation and the exploration of alternative
245 scenarios. By leveraging the computational resources of the cloud, they can execute complex
246 and time-consuming analyses, such as what-if simulations and the evaluation of coordinated
247 strategies, supporting more informed and anticipatory decision-making.

248 Global twins are also the main location where the architecture realizes its full cognitive
249 capabilities. They can exploit more powerful AI models, including full-scale LLMs, to support
250 semantic reasoning, explanation, and human interaction. At this level, LLMs are used to
251 interpret aggregated information, relate observations across subsystems, explain system
252 behavior, and assist in understanding the impact of potential actions. In this way, global
253 cognition complements local cognition, providing deeper and broader reasoning over the
254 entire system.

255 The interaction between local and global twins is mediated by a *shared data space*, which
256 acts as the integration layer of the architecture. This component enables information exchange
257 and coordination without tightly coupling the twins. Instead of direct communication, all
258 components interact through a common information space where structured data, events,
259 summaries, and inferred knowledge are stored and accessed. This design supports decoupling
260 and scalability, as twins do not need to be aware of each other's internal structure.

261 Through this mechanism, local twins continuously generate structured information and
262 propagate it to the data space. Global twins consume this information, reason over the system
263 state, and produce higher-level insights, strategies, or recommendations. These can then
264 be propagated back to the edge, where local twins use them to refine their behavior. This
265 establishes a bidirectional flow of information, combining bottom-up knowledge propagation
266 with top-down coordination.

267 Overall, the proposed architecture defines a distributed organization of intelligence.
268 Instead of concentrating all capabilities in the cloud, DT intelligence is explicitly divided
269 across levels according to their role and constraints. At the edge, DT intelligence is lightweight,
270 fast, and closely tied to the physical system, enabling local autonomy and immediate reactions.
271 In the cloud, DT intelligence is more computationally intensive and supports system-level
272 understanding, coordination, and explanation. These capabilities are complementary rather
273 than redundant.

274 The overall proposed architecture simplifies the engineering of complex DT systems by
275 decomposing the CPS into multiple coordinated twins, improves responsiveness by enabling
276 local decisions, reduces communication through semantic distillation, and supports privacy-
277 aware data sharing. At the same time, it enables advanced reasoning by combining the
278 outputs of multiple twins with simulation and LLM-based interpretation.

4 Conclusion and Future Work

This paper introduced the **Federated Cognitive Digital Twin**, a unified architecture that combines federation and cognition to support DTs in large-scale and distributed CPSs. The proposed approach addresses a key gap in current research. On the one hand, federated DT approaches improve structural decomposition and scalability, but often leave intelligence concentrated in centralized cloud components. On the other hand, cognitive DT approaches enhance semantic reasoning and explainability, but are typically designed as centralized solutions and therefore remain difficult to apply in latency-sensitive and distributed settings.

The FCDDT architecture bridges these two lines of work by explicitly distributing DT capabilities across the edge-to-cloud continuum. In the proposed design, local twins operate close to the physical subsystems and provide real-time monitoring, analysis, first-level decision-making, and lightweight cognitive processing. Global twins operate in the cloud and support system-level reasoning, simulation, explanation, and coordination across multiple subsystems. Their interaction through a shared data space enables loose coupling, scalable integration, and bidirectional knowledge exchange. In this way, the architecture combines local autonomy and responsiveness with global semantic reasoning and coordinated decision support.

Overall, the proposed architecture contributes a coherent way to rethink DT engineering in distributed environments. It supports the decomposition of complex CPSs into multiple coordinated twins, reduces the drawbacks of centralized architectures, and extends DTs beyond data processing toward more meaningful and explainable decision support. For these reasons, we believe that FCDDT represents a promising architectural direction for DTs in scenarios such as smart cities and other complex distributed systems.

As future work, we plan to refine the architecture along both engineering and evaluation dimensions. First, we will further formalize the responsibilities, interfaces, and interaction patterns of local and global twins, including the role of the shared data space and the coordination mechanisms between edge and cloud components. Second, we will investigate model allocation strategies for deciding which cognitive capabilities should remain at the edge and which should be delegated to the cloud, considering trade-offs among latency, bandwidth, computational cost, privacy, and quality of reasoning. Third, we plan to develop a proof-of-concept implementation and assess the architecture in realistic distributed scenarios, with particular attention to responsiveness, resilience, communication overhead, and the usefulness of semantic distillation and LLM-based reasoning for decision support.

In addition, future work should examine how the proposed architecture can support stronger interoperability and governance. This includes studying semantic integration mechanisms across heterogeneous twins, privacy-aware information sharing, and explainability techniques that make system-level recommendations more transparent to human operators. Finally, we aim to investigate tool support and model-driven engineering techniques for simplifying the design, deployment, and evolution of FCDDT-based systems, so that the proposed architecture can be adopted more systematically in practice.

Acknowledgment

This work is supported by the Swedish Agency for Innovation Systems through the project "iSecure: Developing Predictable and Secure IoT for Autonomous Systems" (2023-01899), by the Key Digital Technologies Joint Undertaking through the project "MATISSE: Model-based engineering of digital twins for early verification and validation of industrial systems" (101140216), and by the Clean Energy Transition Partnership through the project "FLEXI:

324 Human-centered AI and digital twin powered energy system integration for flexibility markets"
325 (101069750).

326 — **References** —

- 327 **1** Sailesh Abburu, Arne J. Berre, Michael Jacoby, Dumitru Roman, Ljiljana Stojanovic, and
328 Nenad Stojanovic. Cognitwin – hybrid and cognitive digital twins for the process industry. In
329 *2020 IEEE International Conference on Engineering, Technology and Innovation (ICE/ITMC)*,
330 pages 1–8, 2020. doi:10.1109/ICE/ITMC49519.2020.9198403.
- 331 **2** Alessio Bucaioni, Romina Eramo, Luca Berardinelli, Hugo Bruneliere, Benoit Combemale,
332 Djamel Eddine Khelladi, Vittoriano Muttillio, Andrey Sadovykh, and Manuel Wimmer. Multi-
333 partner project: A model-driven engineering framework for federated digital twins of industrial
334 systems (matisse). In *2025 Design, Automation & Test in Europe Conference (DATE)*, pages
335 1–6. IEEE, 2025.
- 336 **3** Fenglong Cai, Dong Yuan, Zhe Yang, and Lizhen Cui. Edge-llm: A collaborative framework
337 for large language model serving in edge computing. In *2024 IEEE International Conference*
338 *on Web Services (ICWS)*, pages 799–809, 2024. doi:10.1109/ICWS62655.2024.00099.
- 339 **4** Alessandra De Benedictis, Nicola Mazzocca, Alessandra Somma, and Carmine Strigaro.
340 Digital twins in healthcare: An architectural proposal and its application in a social distancing
341 case study. *IEEE Journal of Biomedical and Health Informatics*, 27(10):5143–5154, 2023.
342 doi:10.1109/JBHI.2022.3205506.
- 343 **5** Kubra Duran, Lal Verda Cakir, Yagmur Yigit, Khayal Huseynov, Sushmitha Ram Kusu,
344 Mehmet Ali Ertürk, and Berk Canberk. Toward digital twin-as-a-service (dtaas) platforms: A
345 survey on architecture, design requirements, and performance metrics. *IEEE Communications*
346 *Surveys & Tutorials*, 28:1845–1878, 2026. doi:10.1109/COMST.2025.3635582.
- 347 **6** Enxhi Ferko, Alessio Bucaioni, and Moris Behnam. Supporting technical adaptation and
348 implementation of digital twins in manufacturing. In *International Conference on Information*
349 *Technology-New Generations*, pages 181–189. Springer, 2012.
- 350 **7** Enxhi Ferko, Alessio Bucaioni, and Moris Behnam. Architecting digital twins. *IEEE Access*,
351 10:50335–50350, 2022. doi:10.1109/ACCESS.2022.3172964.
- 352 **8** Enxhi Ferko, Alessio Bucaioni, Patrizio Pelliccione, and Moris Behnam. Analysing interoper-
353 ability in digital twin software architectures for manufacturing. In *European conference on*
354 *software architecture*, pages 170–188. Springer, 2023.
- 355 **9** Enxhi Ferko, Alessio Bucaioni, Patrizio Pelliccione, and Moris Behnam. Standardisation in
356 digital twin architectures in manufacturing. In *2023 IEEE 20th International Conference on*
357 *Software Architecture (ICSA)*, pages 70–81. IEEE, 2023.
- 358 **10** Othmane Friha, Mohamed Amine Ferrag, Burak Kantarci, Burak Cakmak, Arda Ozgun, and
359 Nassira Ghoulmi-Zine. Llm-based edge intelligence: A comprehensive survey on architectures,
360 applications, security and trustworthiness. *IEEE Open Journal of the Communications Society*,
361 5:5799–5856, 2024. doi:10.1109/OJCOMS.2024.3456549.
- 362 **11** Zixu Hao, Huiqiang Jiang, Shiqi Jiang, Ju Ren, and Ting Cao. Hybrid slm and llm for edge-
363 cloud collaborative inference. In *Proceedings of the Workshop on Edge and Mobile Foundation*
364 *Models, EdgeFM '24*, page 36–41, New York, NY, USA, 2024. Association for Computing
365 Machinery. doi:10.1145/3662006.3662067.
- 366 **12** Yangyang Liu, Tang Ji, Xiangyu Guo, Xun Xu, and Jan Polzer. A survey of cognitive digital
367 twin and the potential use of llms. *Manufacturing Letters*, 44:1242–1253, 2025. 53rd SME
368 North American Manufacturing Research Conference (NAMRC 53). doi:10.1016/j.mfglet.
369 2025.06.144.
- 370 **13** Andrea Manzi, Raul Bardaji, Ivan Rodero, Germán Moltó, Sandro Fiore, Isabel Campos,
371 Donatello Elia, Francesco Sarandrea, A. Paul Millar, Daniele Spiga, Matteo Bunino, Gabriele
372 Accarino, Lorenzo Asprea, Samuel Bernardo, Miguel Caballer, Charis Chatzikyriakou, Diego
373 Ciangottini, Michele Claus, Andrea Cristofori, Davide Donno, Emanuele Donno, Iacopo

- 374 Ferrario, Massimiliano Fronza, Alexander Jacob, Javad Komijani, Marina Krstic Marinkovic,
375 Federica Legger, Ivan Palomo, Estíbaliz Parceró, Rakesh Sarma, Gaurav Sinha Ray, Sara
376 Vallero, and Juraj Zvolensky. intertwin: Advancing scientific digital twins through ai, federated
377 computing and data. *Future Generation Computer Systems*, 179:108312, 2026. doi:10.1016/
378 j.future.2025.108312.
- 379 14 Judith Michael, Loek Cleophas, Steffen Zschaler, Tony Clark, Benoit Combemale, Thomas
380 Godfrey, Djamel Eddine Khelladi, Vinay Kulkarni, Daniel Lehner, Bernhard Rumpe, Manuel
381 Wimmer, Andreas Wortmann, Shaukat Ali, Balbir Barn, Ion Barosan, Nelly Bencomo, Francis
382 Bordeleau, Georg Grossmann, Gabor Karsai, Oliver Kopp, Bernhard Mitschang, Paula
383 Muñoz Ariza, Alfonso Pierantonio, Fiona A. C. Polack, Matthias Riebisch, Holger Schling-
384 gloff, Markus Stumptner, Antonio Vallecillo, Mark van den Brand, and Hans Vangheluwe.
385 Model-driven engineering for digital twins: Opportunities and challenges. *Systems Engineering*,
386 28(5):659–670, 2025. doi:10.1002/sys.21815.
- 387 15 Roberto Minerva, Gyu Myoung Lee, and Noël Crespi. Digital twin in the iot context: A
388 survey on technical features, scenarios, and architectural models. *Proceedings of the IEEE*,
389 108(10):1785–1824, 2020. doi:10.1109/JPROC.2020.2998530.
- 390 16 Alexios Papacharalampopoulos, Dionysios Christopoulos, Olga Maria Karagianni, and
391 Panagiotis Stavropoulos. Federation in digital twins and knowledge transfer: Modeling
392 limitations and enhancement. *Machines*, 12(10), 2024. doi:10.3390/machines12100701.
- 393 17 Renuka Prasad Pasupulati and Jordan Shropshire. Analysis of centralized and decentralized
394 cloud architectures. In *SoutheastCon 2016*, pages 1–7, 2016. doi:10.1109/SECON.2016.
395 7506680.
- 396 18 Kunal Rao, Giuseppe Coviello, Priscilla Benedetti, Ciro Giuseppe De Vita, Gennaro Mellone,
397 and Srimat Chakradhar. Eco-llm: Llm-based edge cloud optimization. In *Proceedings of
398 the 2024 Workshop on AI For Systems*, AI4Sys '24, page 7–12. Association for Computing
399 Machinery, 2024. doi:10.1145/3660605.3660941.
- 400 19 Daniel Rosendo, Alexandru Costan, Patrick Valduriez, and Gabriel Antoniu. Distributed
401 intelligence on the edge-to-cloud continuum: A systematic literature review. *Journal of Parallel
402 and Distributed Computing*, 166:71–94, 2022. doi:10.1016/j.jpdc.2022.04.004.
- 403 20 Mahadev Satyanarayanan. The emergence of edge computing. *Computer*, 50(1):30–39, 2017.
404 doi:10.1109/MC.2017.9.
- 405 21 Weisong Shi, Jie Cao, Quan Zhang, Youhuizi Li, and Lanyu Xu. Edge computing: Vision and
406 challenges. *IEEE Internet of Things Journal*, 3(5):637–646, 2016. doi:10.1109/JIOT.2016.
407 2579198.
- 408 22 Inderjeet Singh, Eleonore Vissol-Gaudin, Andikan Otung, and Motoyoshi Sekiya. Learning
409 to collaborate: An orchestrated-decentralized framework for peer-to-peer llm federation.
410 *Proceedings of the AAAI Conference on Artificial Intelligence*, 40(30):25472–25480, Mar. 2026.
411 doi:10.1609/aaai.v40i30.39742.
- 412 23 Alessandra Somma, Domenico Amalfitano, Alessio Bucaioni, and Alessandra De Benedictis.
413 A model-driven approach for engineering mobility digital twins: The bologna case study.
414 *Information and Software Technology*, page 107863, 2025.
- 415 24 Alessandra Somma, Domenico Amalfitano, Alessandra De Benedictis, and Patrizio Pelliccione.
416 Twinarch: A digital twin reference architecture. *Journal of Systems and Software*, 231:112613,
417 2026. doi:/10.1016/j.jss.2025.112613.
- 418 25 Fei Tao, Bin Xiao, Qinglin Qi, Jiangfeng Cheng, and Ping Ji. Digital twin modeling. *Journal
419 of Manufacturing Systems*, 64:372–389, 2022. doi:10.1016/j.jmsy.2022.06.015.
- 420 26 Christian Vergara, Rami Bahsoon, Georgios Theodoropoulos, Wendy Yanez, and Nikos
421 Tziritas. Federated digital twin. In *2023 IEEE/ACM 27th International Symposium on
422 Distributed Simulation and Real Time Applications (DS-RT)*, pages 115–116, 2023. doi:
423 10.1109/DS-RT58998.2023.00027.
- 424 27 Christian Roberto Vergara, Georgios Theodoropoulos, Rami Bahsoon, Wendy Yanez, and
425 Nikos Tziritas. Federated digital twins as an enabling technology for collaborative decision-

- 426 making. In *Proceedings of the 38th ACM SIGSIM Conference on Principles of Advanced*
427 *Discrete Simulation*, SIGSIM-PADS '24, page 67–68. Association for Computing Machinery,
428 2024. doi:10.1145/3615979.3662152.
- 429 **28** Christian Vergara-Marcillo, Rami Bahsoon, Nikos Tziritas, and Georgios Theodoropoulos. A
430 connectionist approach to federated digital twins. In *Computational Science – ICCS 2025*,
431 pages 60–74, Cham, 2025. Springer Nature Switzerland.
- 432 **29** Yuchen Xia, Nasser Jazdi, and Michael Weyrich. An architecture for integrating large
433 language models with digital twins and automation systems. In *2025 IEEE 30th International*
434 *Conference on Emerging Technologies and Factory Automation (ETFA)*, pages 1–8, 2025.
435 doi:10.1109/ETFA65518.2025.11205636.
- 436 **30** Nan Zhang, Rami Bahsoon, and Georgios Theodoropoulos. Towards engineering cognitive
437 digital twins with self-awareness. In *2020 IEEE International Conference on Systems, Man,*
438 *and Cybernetics (SMC)*, pages 3891–3891, 2020. doi:10.1109/SMC42975.2020.9283357.
- 439 **31** Xiaochen Zheng, Jinzhi Lu, and Dimitris Kiritsis. The emergence of cognitive digital twin:
440 vision, challenges and opportunities. *International Journal of Production Research*, 60(24):7610–
441 7632, 2022. doi:10.1080/00207543.2021.2014591.
- 442 **32** Fang Zhu, Jiayuan Chen, Junjie Wen, Yuye Yang, Changyan Yi, Yun Tie, Peng Zhang, Jun
443 Cai, Dusit Niyato, and Mohsen Guizani. From data mirror to smart copilot: A survey on
444 nextg semantic communication for propelling digital twin world into cognitive stage. *IEEE*
445 *Communications Surveys & Tutorials*, 28:4915–4947, 2026. doi:10.1109/COMST.2026.3665395.