

Edge Intelligence: the Confluence of Edge Computing and Artificial Intelligence

Hailiang Zhao
hliangzhao@zju.edu.cn

College of Computer Science and Technology, Zhejiang University

November 17, 2019

Outline

1 Introduction

- 5G, edge, and AI
- Relations between Edge Computing and AI
- Birth of Edge Intelligence

Outline

1 Introduction

- 5G, edge, and AI
- Relations between Edge Computing and AI
- Birth of Edge Intelligence

2 Research Roadmap of Edge Intelligence

- Roadmap overview
- Quality of Experience
- Intelligence-enabled Edge Computing
- Artificial Intelligence on Edge

Outline

1 Introduction

- 5G, edge, and AI
- Relations between Edge Computing and AI
- Birth of Edge Intelligence

2 Research Roadmap of Edge Intelligence

- Roadmap overview
- Quality of Experience
- Intelligence-enabled Edge Computing
- Artificial Intelligence on Edge

3 AI for Edge

- State of the Art
- Grand Challenges

Outline

1 Introduction

- 5G, edge, and AI
- Relations between Edge Computing and AI
- Birth of Edge Intelligence

2 Research Roadmap of Edge Intelligence

- Roadmap overview
- Quality of Experience
- Intelligence-enabled Edge Computing
- Artificial Intelligence on Edge

3 AI for Edge

- State of the Art
- Grand Challenges

4 AI on Edge

- State of the Art
- Grand Challenges

Outline

1 Introduction

- 5G, edge, and AI
- Relations between Edge Computing and AI
- Birth of Edge Intelligence

2 Research Roadmap of Edge Intelligence

- Roadmap overview
- Quality of Experience
- Intelligence-enabled Edge Computing
- Artificial Intelligence on Edge

3 AI for Edge

- State of the Art
- Grand Challenges

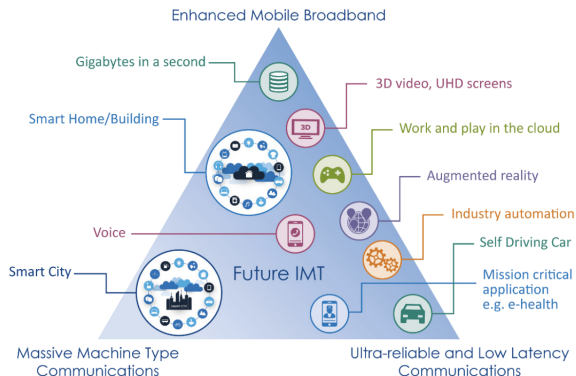
4 AI on Edge

- State of the Art
- Grand Challenges

5G is coming!

What 5G brings to us

- 1 enhanced Mobile BroadBand (**eMBB**)
- 2 Ultra-Reliable Low Latency Communications (**URLLC**)
- 3 massive Machine Type Communications (**mMTC**)



Processing data nearby¹

Why **edge**?

- ① **explosion** of data generated by mobile and IoT devices
- ② **oppressive** network congestion in backbone
- ③ ...

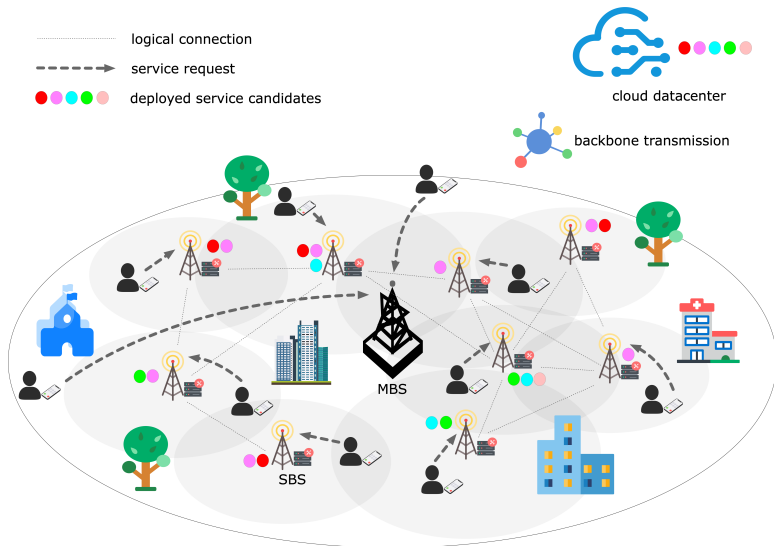
Multi-access Edge Computing (MEC)

- ① communication/computation/caching/control at the edge directly
- ② provide services
- ③ perform computations
- ④ manage resources

MEC avoids unnecessary communication latency and enabling faster responses for end users.

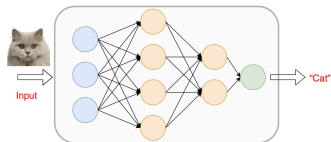
¹Z. Zhou et al. "Edge Intelligence: Paving the Last Mile of Artificial Intelligence With Edge Computing". In: *Proceedings of the IEEE* 107.8 (2019), pp. 1738–1762.

A typical pre-5G HetNet

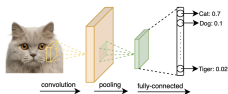


What about Artificial Intelligence?

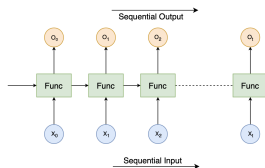
- ① powerfull in **big data processing & insights extracting**
- ② DNNs: powerfull **knowledge representation**
- ③ Typical structures of DNNs
 - ① Multilayer Perceptrons (MLP)
 - ② Convolutional Neural Network (CNN) (AlexNet \rightarrow VGG-16 \rightarrow GoogleNet \rightarrow ResNet)
 - ③ Recurrent Neural Network (RNN) (RNN \rightarrow LSTM)
- ④ Popular DNN models
 - ① Generative Adversarial Network (GAN)
 - ② Deep Reinforcement Learning (DRL)



(a) Multilayer Perceptrons



(b) Convolution Neural Network



(c) Recurrent Neural Network

Can they integrate with each other?

- ① AI provides Edge Computing with **methods and technologies**
 - ① Complicated resource allocation problems need to solve
 - ② Huge volumes of data need to analysis
 - ③ AI can help in **model formulation & optimization**
- ② Edge Computing provides AI with **scenarios and platforms**
 - ① More and more data is created by widespread and **geographically distributed** mobile and IoT devices
 - ② Many more applicaiton **scenarios** (intelligent networked vehicles, autonomous driving, smart hone, smart city, ...)
 - ③ **Hardware acceleration** on resource-limited IoT devices

Their integration leads to the birth of

Edge Intelligence (a.k.a. Edge AI)

Edge Intelligence: our definition

Edge Intelligence

We divide it into **AI for edge** and **AI on edge**.

① **AI for edge**

- ① provide a better solution to the constrained optimization problems
- ② AI is used for energizing edge with more intelligence and optimality
- ③ **Intelligence-enabled Edge Computing (IEC)**

② **AI on edge**

- ① carry out the entire process of AI models on edge
- ② run model training and inference with device-edge-cloud synergy
- ③ **Artificial Intelligence on Edge (AIE)**

Outline

1 Introduction

- 5G, edge, and AI
- Relations between Edge Computing and AI
- Birth of Edge Intelligence

2 Research Roadmap of Edge Intelligence

- Roadmap overview
- Quality of Experience
- Intelligence-enabled Edge Computing
- Artificial Intelligence on Edge

3 AI for Edge

- State of the Art
- Grand Challenges

4 AI on Edge

- State of the Art
- Grand Challenges

Roadmap Overview



QoE: indicators

① performance

- ① AI for edge: problem-dependent
- ② AI on edge: training loss, inference loss

② cost

- ① computation cost (CPU time, CPU frequency)
- ② communication cost (transmit power, frequency band, access time)
- ③ energy consumption (battery capacity)

③ privacy (security)

- ① leads to the birth of **Federated Learning**

④ efficiency

- ① excellent performance with low overhead

⑤ reliability

- ① robustness
- ② handle with failure

AI for edge: a recapitulation

1 Service

- 1 optimize computation offloading via DQN²³

2 Content

- 1 service placement via MAB⁴
- 2 service deployment via DRL⁵

3 Topology

- 1 optimize UAVs via Multi-agent Learning⁶
- 2 learning-driven communication⁷

²X. Chen et al. “Optimized Computation Offloading Performance in Virtual Edge Computing Systems Via Deep Reinforcement Learning”. In: *IEEE Internet of Things Journal* 6.3 (2019), pp. 4005–4018.

³M. Min et al. “Learning-Based Computation Offloading for IoT Devices With Energy Harvesting”. In: *IEEE Transactions on Vehicular Technology* 68.2 (2019), pp. 1930–1941.

⁴L. Chen et al. “Spatio-Temporal Edge Service Placement: A Bandit Learning Approach”. In: *IEEE Transactions on Wireless Communications* 17.12 (2018), pp. 8388–8401.

⁵Y. Chen et al. “Data-Intensive Application Deployment at Edge: A Deep Reinforcement Learning Approach”. In: *2019 IEEE International Conference on Web Services (ICWS)*. 2019, pp. 355–359.

⁶J. Xu, Y. Zeng, and R. Zhang. “UAV-Enabled Wireless Power Transfer: Trajectory Design and Energy Optimization”. In: *IEEE Transactions on Wireless Communications* 17.8 (2018), pp. 5092–5106.

⁷M. Chen et al. “Artificial Neural Networks-Based Machine Learning for Wireless Networks: A Tutorial”. In: *IEEE Communications Surveys Tutorials* (2019), pp. 1–33.

AI on edge: a recapitulation

① model adaptation (too many of them)

- ① model compression, conditional computation, algorithm asynchronization, thoroughly decentralization, ...

② framework design

- ① model training: Federated Learning on edge⁸, knowledge distillation-based methods⁹
- ② model inference: model splitting/partitioning (Edgent)¹⁰

③ processor acceleration¹¹

- ① design special instruction sets
- ② design high parallel computing paradigms
- ③ move computation closer to memory

⁸Kai Yang et al. "Federated Learning via Over-the-Air Computation". In: *CoRR abs/1812.11750* (2018). arXiv: 1812.11750.

⁹Jin-Hyun Ahn, Osvaldo Simeone, and Joonhyuk Kang. "Wireless Federated Distillation for Distributed Edge Learning with Heterogeneous Data". In: *ArXiv abs/1907.02745* (2019).

¹⁰En Li, Zhi Zhou, and Xu Chen. "Edge Intelligence: On-Demand Deep Learning Model Co-Inference with Device-Edge Synergy". In: *Proceedings of the 2018 Workshop on Mobile Edge Communications, MECOMM@SIGCOMM 2018, Budapest, Hungary, August 20, 2018*. 2018, pp. 31–36.

¹¹V. Sze et al. "Efficient Processing of Deep Neural Networks: A Tutorial and Survey". In: *Proceedings of the IEEE* 105.12 (2017), pp. 2295–2329.

Outline

1 Introduction

- 5G, edge, and AI
- Relations between Edge Computing and AI
- Birth of Edge Intelligence

2 Research Roadmap of Edge Intelligence

- Roadmap overview
- Quality of Experience
- Intelligence-enabled Edge Computing
- Artificial Intelligence on Edge

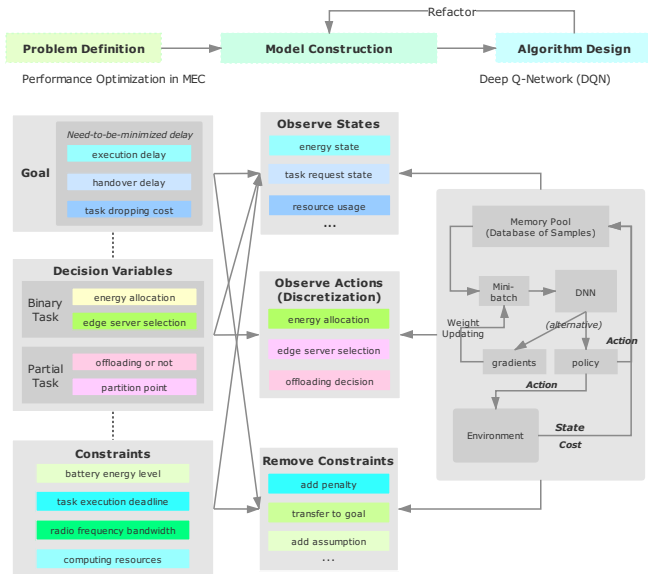
3 AI for Edge

- State of the Art
- Grand Challenges

4 AI on Edge

- State of the Art
- Grand Challenges

Utilize DQN for performance optimization



Grand challenges

① model establishment

- ① unrestrained searching space
- ② state/action set cannot be infinite

② algorithm deployment

- ① cannot obtain analytic (approximate) optimal solution
- ② too many iterations → hard to deploy in an online manner
- ③ who undertake the responsibility?

③ balance between optimality and efficiency

Outline

1 Introduction

- 5G, edge, and AI
- Relations between Edge Computing and AI
- Birth of Edge Intelligence

2 Research Roadmap of Edge Intelligence

- Roadmap overview
- Quality of Experience
- Intelligence-enabled Edge Computing
- Artificial Intelligence on Edge

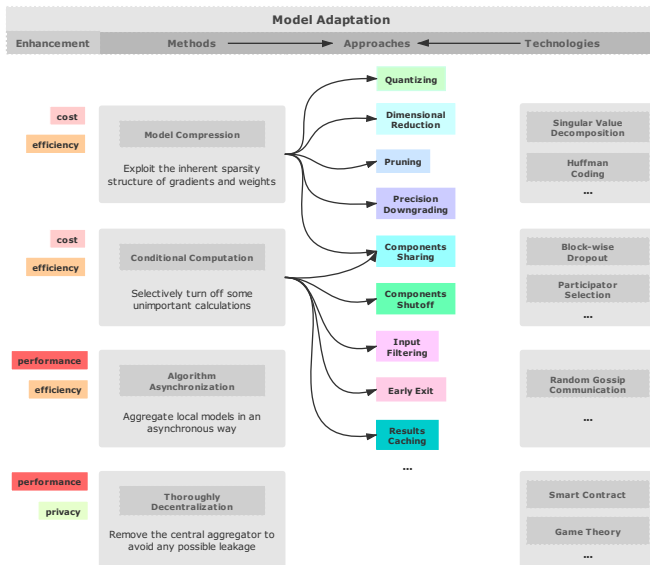
3 AI for Edge

- State of the Art
- Grand Challenges

4 AI on Edge

- State of the Art
- Grand Challenges

Model Adaptation: a classification



Grand challenges

① data availability

- ① where to find **usable** data?
- ② incentive mechanisms
- ③ obvious bias from distributed end users (non i.i.d.)

② model selection

- ① select befitting threshold of learning accuracy & scale of models
- ② select probe training frameworks and accelerator architectures

③ coordination mechanism

- ① same method achieves different results
- ② compatibility and coordination (cloud-edge-device synergy)
- ③ establish a unified API interface?