

Deep Learning for **Wireless Networks**

-- Which Model to Use?

Jun ZHANG



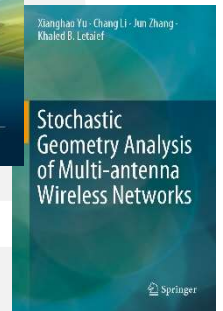
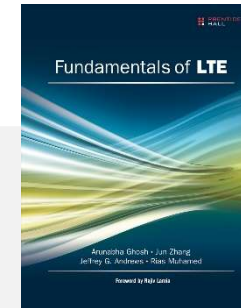
THE HONG KONG
POLYTECHNIC UNIVERSITY
香港理工大學

My Research Interests



Wireless Communications and Networking

5G and beyond,
Internet of Things



Edge AI, Mobile AI

Mobile edge computing,
Edge learning,
Federated learning



Big Data Analytics

Distributed machine learning,
Cloud computing

Interested in 5G?

- **Sparse and Low-Rank Optimization for Dense Wireless Networks**

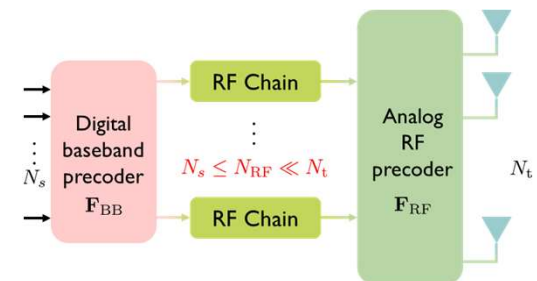
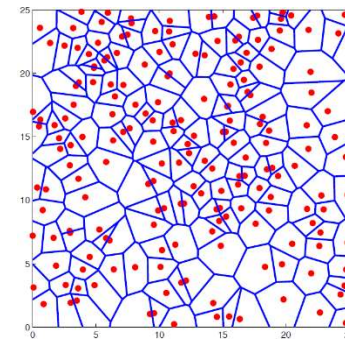
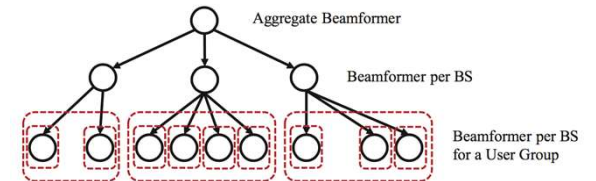
- Tutorial at IEEE GLOBECOM 2017

- **Tractable Analysis of Large-scale Multi-antenna Wireless Networks via Stochastic Geometry**

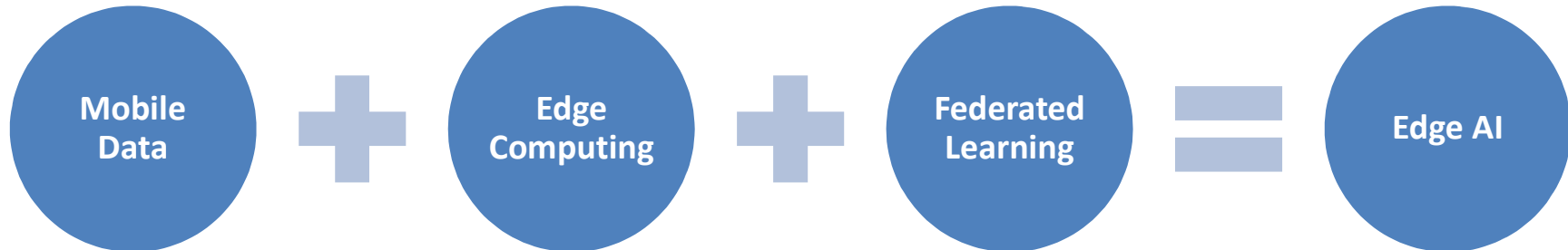
- Tutorial at WiOpt 2018

- **Hybrid Beamforming for 5G Millimeter Wave Systems**

- Tutorial at IEEE GLOBECOM 2018

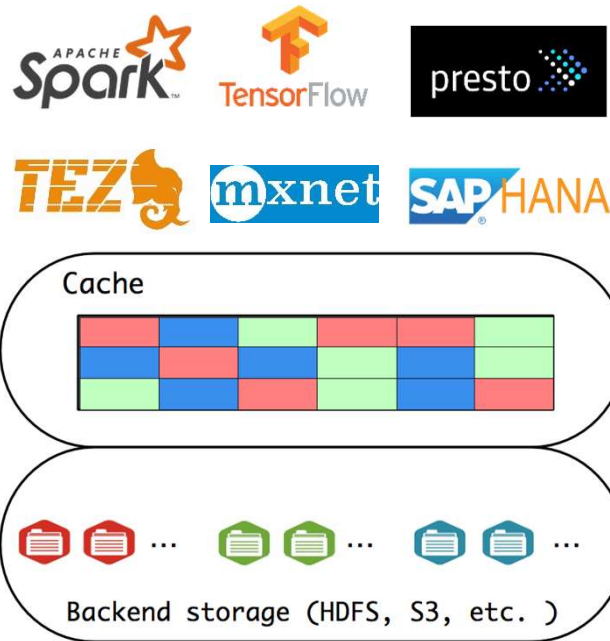


Edge AI



- Y. Mao, C. You, **J. Zhang**, K. Huang, and K. B. Letaief, “**A survey on mobile edge computing: The communication perspective**,” *IEEE Commun. Surveys Tuts.*, vol. 19, no. 4, pp. 2322-2358, 4th Quart. 2017.
- G. Zhu, D. Liu, Y. Du, C. You, **J. Zhang**, and K. Huang, “**Towards an intelligent edge: Wireless communication meets machine learning**,” *submitted to IEEE Communications Magazine*.
- **J. Zhang**, and K. B. Letaief, “**Mobile Edge Intelligence and Computing for the Internet of Vehicles**,” *submitted to Proc. IEEE*.

Big Data Analytics – Cache Management



Cache is limited – What to cache?

- Y.Yu, W.Wang, **J. Zhang**, and K. B. Letaief, “LRC: Dependency-aware cache management in data analytics clusters,” in *Proc. IEEE INFOCOM 2017*. (Acceptance Rate: 20.93%)

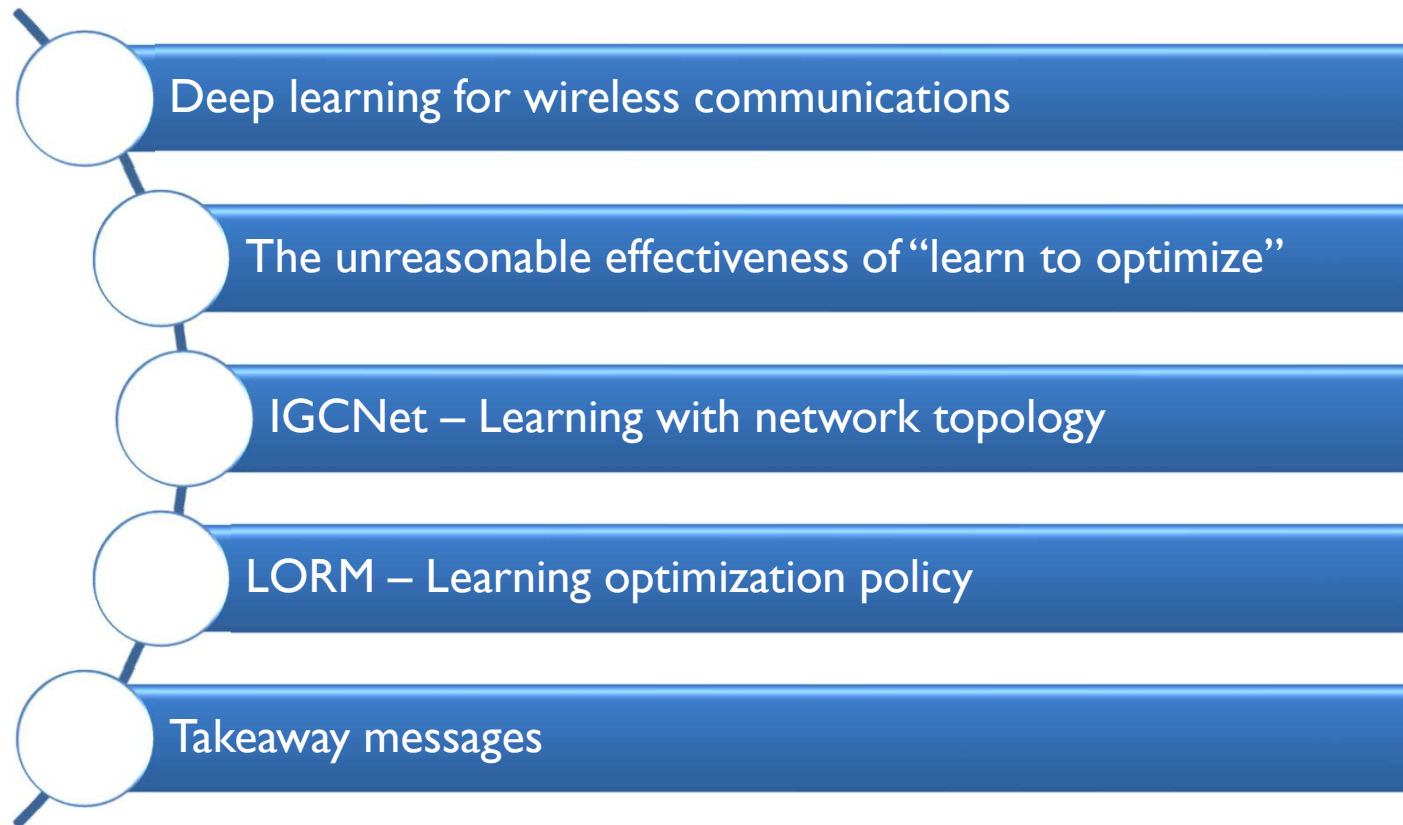
Cache is shared resource – How to fairly share cache among users?

- Y.Yu, W.Wang, **J. Zhang**, Q. Weng, and K. B. Letaief, “OpuS: Fair and efficient cache sharing for in-memory data analytics,” in *ICDCS 2018*. (Acceptance Rate: 20%)

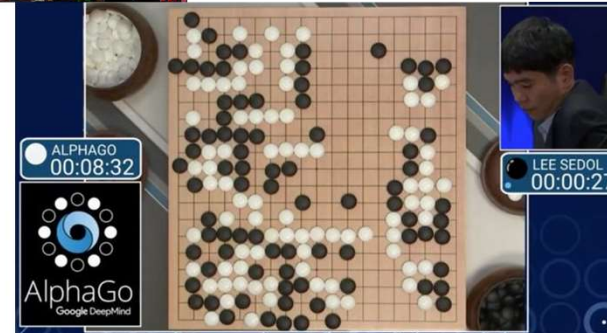
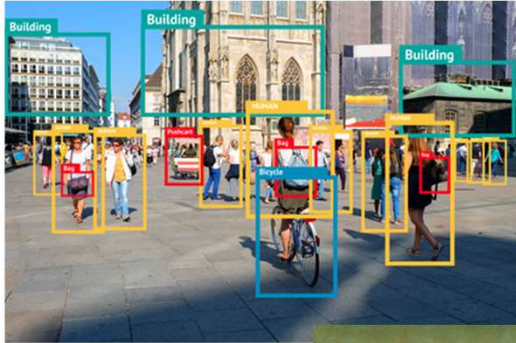
Distributed caches – How to balance the load?

- Y.Yu, R. Huang, W.Wang, **J. Zhang**, and K. B. Letaief, “SP-Cache: Load-balanced, Redundancy-free Cluster Caching with Selective Partition,” in *SC 2018*. (Acceptance Rate: 19%)
- Y.Yu, W.Wang, **J. Zhang**, and K. B. Letaief, “LACS: Load-aware cache sharing with isolation guarantee,” *ICDCS 2019*. (Acceptance Rate: 19.6%)

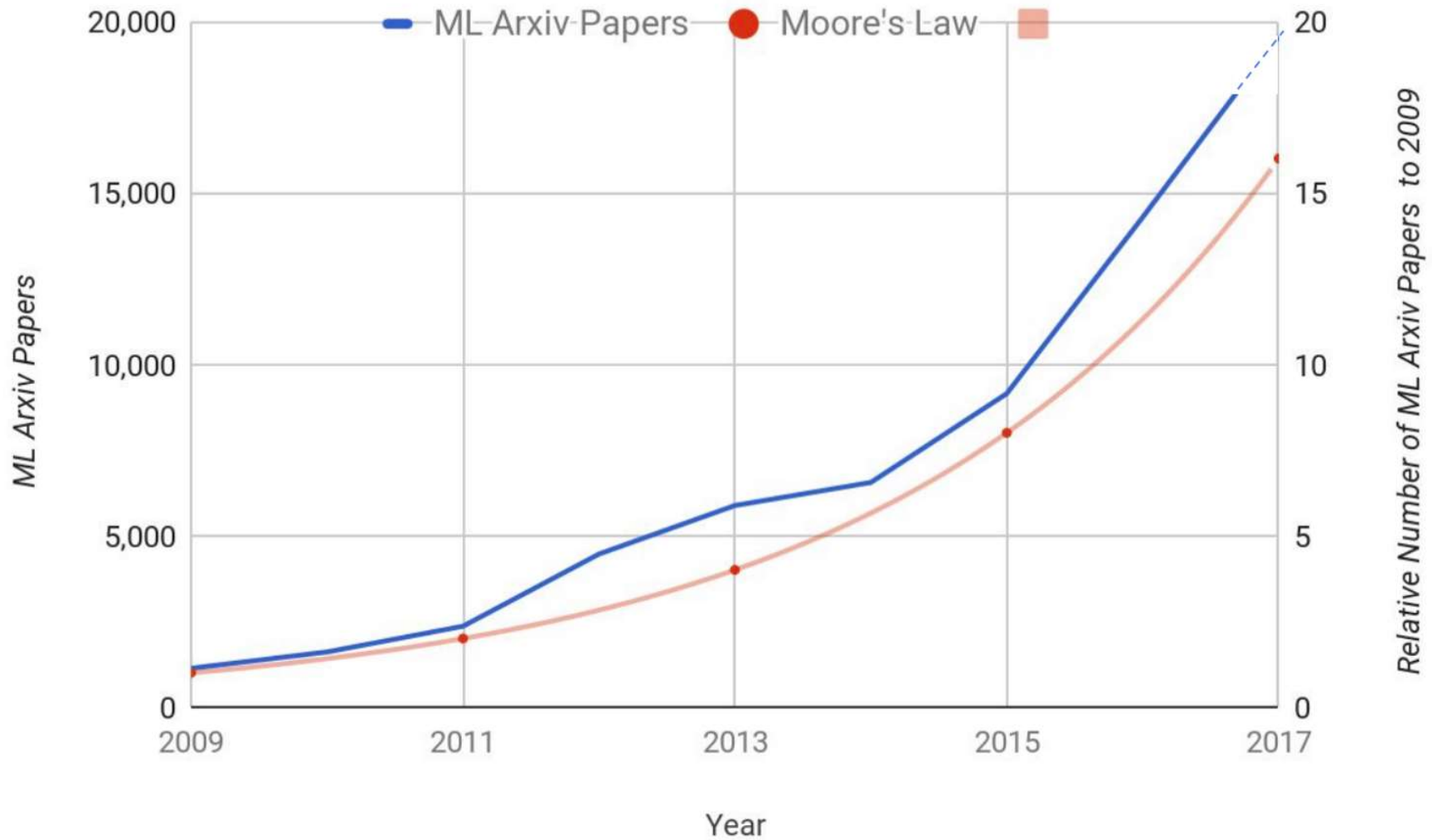
Outline



Successes of Deep Learning



ML Arxiv Papers per Year



Deep Learning: Alchemy or Science?



Deep Learning in Wireless Communications

Deep Learning for Wireless Physical Layer: Opportunities and Challenges

Nov 2017

Tianqi Wang¹, Chao-Kai Wen², Hanqing Wang¹, Feifei Gao³, Tao Jiang⁴, Shi Jin^{1,*}

IEEE TRANSACTIONS ON COGNITIVE COMMUNICATIONS AND NETWORKING, VOL. 3, NO. 4, DECEMBER 2017

563

An Introduction to Deep Learning for the Physical Layer

Dec 2017

Timothy O'Shea¹, *Senior Member, IEEE*, and Jakob Hoydis, *Member, IEEE*

The Roadmap to 6G: AI Empowered Wireless Networks

Aug 2019

Khaled B. Letaief, Wei Chen, Yuanming Shi, Jun Zhang, and Ying-Jun Angela Zhang

Wireless Networks Design in the Era of Deep Learning: Model-Based, AI-Based, or Both?

To appear

Alessio Zappone, *Senior Member, IEEE*, Marco Di Renzo, *Senior Member, IEEE*, M rouane Debbah, *Fellow, IEEE*
(Invited Paper)

MODEL-DRIVEN DEEP LEARNING FOR PHYSICAL LAYER COMMUNICATIONS

To appear

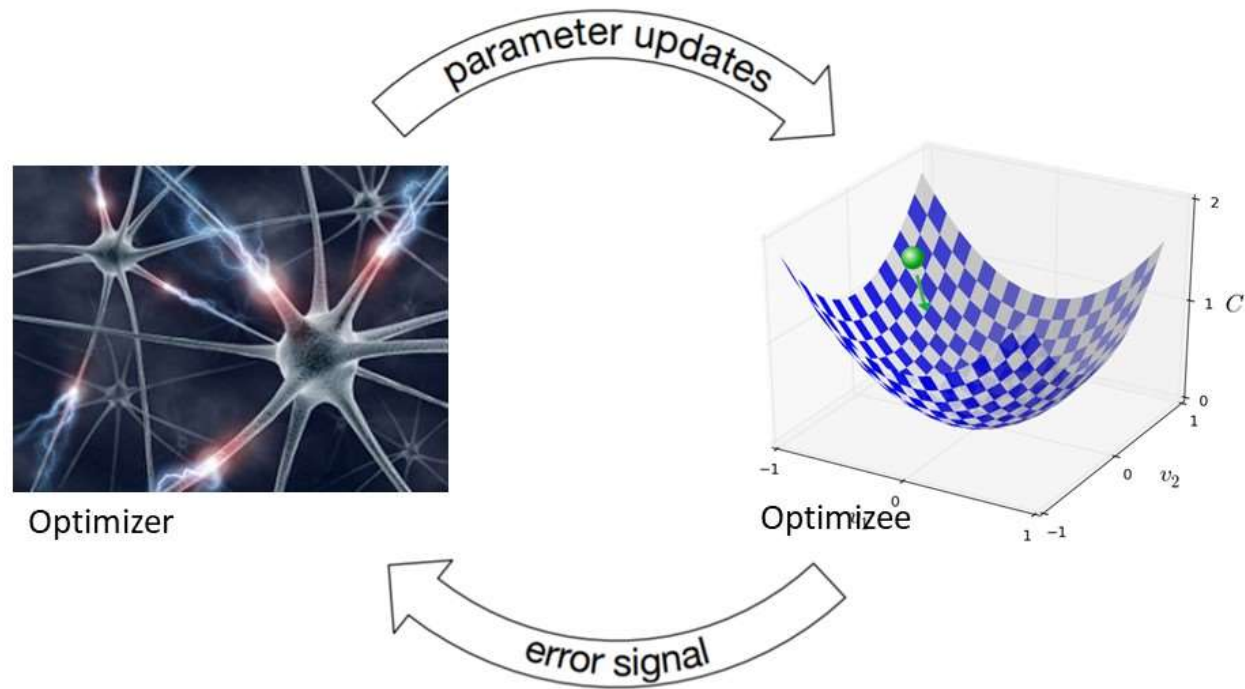
Hengtao He, Shi Jin, Chao-Kai Wen, Feifei Gao, Geoffrey Ye Li, and Zongben Xu

When can deep learning help?

- **When it is difficult to obtain good models**
 - Channel estimation in mmWave systems
 - Traffic prediction, user mobility
- **When there is a lack of design methodology, but abundant data**
 - Channel coding with feedback
 - Joint source-channel coding
- **When conventional methods work, but too complex**
 - Learn to optimize for resource management

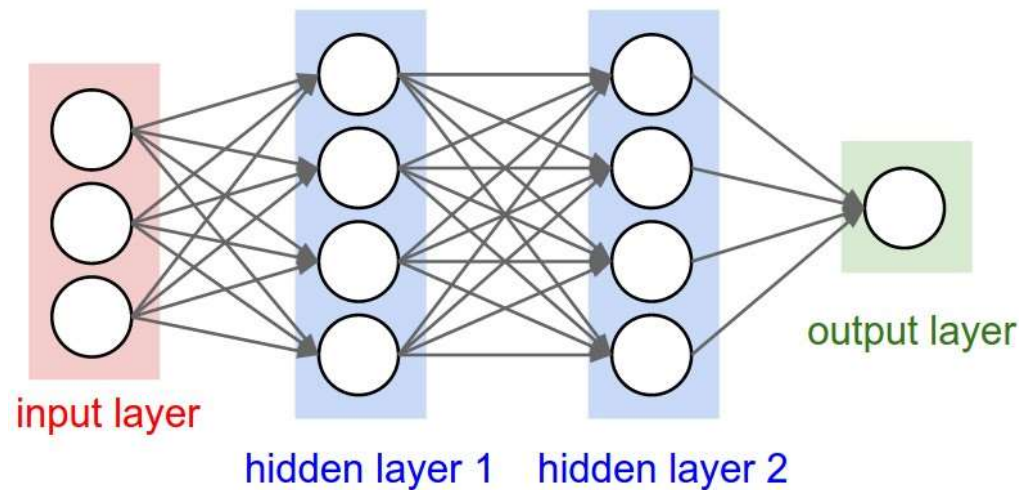
The focus of this talk

Unreasonable effectiveness of “learn to optimize”



Supervised Learning with Neural Networks

$$\underset{\Theta}{\text{minimize}} \quad \ell(f(\Theta, \mathbf{X}))$$



k-th layer's output

$$\mathbf{g}^k = \text{Relu}(\mathbf{W}^k \mathbf{g}^{k-1} + \mathbf{b}^k),$$

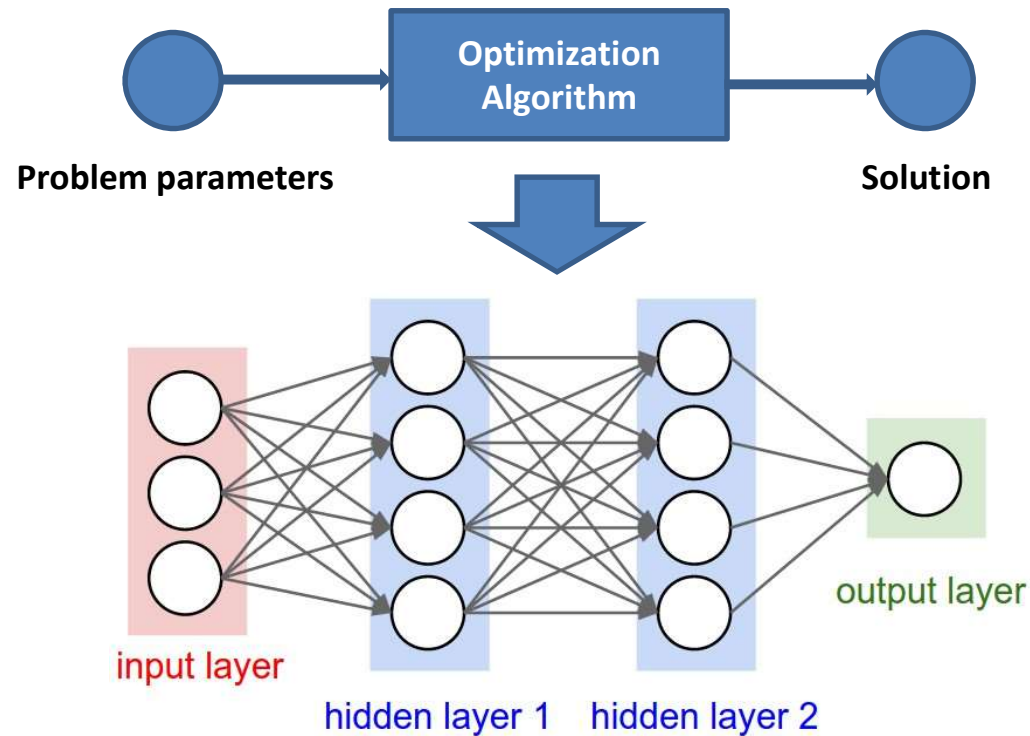
\mathbf{g}^k is the output of k-th layer

$\mathbf{W}^k, \mathbf{b}^k$ the learned parameters

$$\text{Relu}(\cdot) = \max(\cdot, 0)$$

Learning to Optimize (L2O)

- Use machine learning techniques to find near-optimal solutions at affordable cost



Recent Interests From ML Community

Learning to Search in Branch-and-Bound Algorithms*

NIPS 2014

He He Hal Daumé III
Department of Computer Science
University of Maryland
College Park, MD 20740
{hhe, hal}@cs.umd.edu

Jason Eisner
Department of Computer Science
Johns Hopkins University
Baltimore, MD 21218
jason@cs.jhu.edu

Learning Combinatorial Optimization Algorithms over Graphs

NIPS 2017

Hanjun Dai^{†*}, Elias B. Khalil^{†*}, Yuyu Zhang[†], Bistra Dilkina[†], Le Song^{†§}
[†] College of Computing, Georgia Institute of Technology
[§] Ant Financial
{hanjun.dai, elias.khalil, yuyu.zhang, bdilkina, lsong}@cc.gatech.edu

Learning to Branch

ICML 2018

Maria-Florina Balcan¹ Travis Dick¹ Tuomas Sandholm¹ Ellen Vitercik¹

Machine Learning for Combinatorial Optimization:
a Methodological Tour d'Horizon*

A recent survey

Yoshua Bengio^{2,3}, Andrea Lodi^{1,3}, and Antoine Prouvost^{1,3}

Optimization Problems are Ubiquitous in Wireless Networks



Input

- CSI
- QoS requirement
- Resource constraints
- etc

Output

- Resource allocation
- Detection results
- Estimates
- etc

Challenges

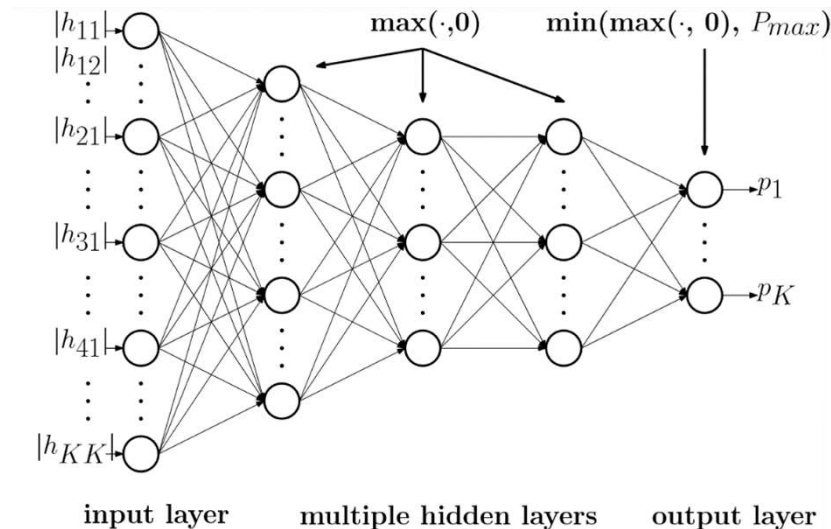
- Large problem size
- Non-convexity
- Real-time execution
- Uncertainty in parameters
- etc

An attempt in wireless networks

$$\max_{p_1, \dots, p_K} \sum_{k=1}^K \alpha_k \log \left(1 + \frac{|h_{kk}|^2 p_k}{\sum_{j \neq k} |h_{kj}|^2 p_j + \sigma_k^2} \right)$$

s.t. $0 \leq p_k \leq P_{\max}, \forall k = 1, 2, \dots, K,$

- WMMSE is a classic algorithm, good performance but slow
- To speed up, use **multi-layer perceptron (MLP)** to approximate the output of WMMSE [Sun I8TSP]



In theory, it works

Theorem 2 Suppose that WMMSE is randomly initialized with $(v_k^0)^2 \leq P_{\max}$, $\sum_{i=1}^K v(h)_i^0 \geq V_{\min}$, and it is executed for T iterations. Define the following set of ‘admissible’ channel realizations

$$\mathcal{H} := \left\{ h \mid H_{\min} \leq |h_{jk}| \leq H_{\max}, \forall j, k, \sum_{i=1}^K v(h)_i^t \geq V_{\min}, \forall t \right\}.$$

Given $\epsilon > 0$, there exists a neural network with $h \in \mathbb{R}^{K^2}$ and $v^0 \in \mathbb{R}_+^K$ as input and $NET(h, v^0) \in \mathbb{R}_+^K$ as output, with the following number of layers

$$O\left(T^2 \log\left(\max\left(K, P_{\max}, H_{\max}, \frac{1}{\sigma}, \frac{1}{H_{\min}}, \frac{1}{P_{\min}}\right)\right) + T \log\left(\frac{1}{\epsilon}\right)\right)$$

and the following number of ReLUs and binary units

$$O\left(T^2 K^2 \log\left(\max\left(K, P_{\max}, H_{\max}, \frac{1}{\sigma}, \frac{1}{H_{\min}}, \frac{1}{P_{\min}}\right)\right) + TK^2 \log\left(\frac{1}{\epsilon}\right)\right),$$

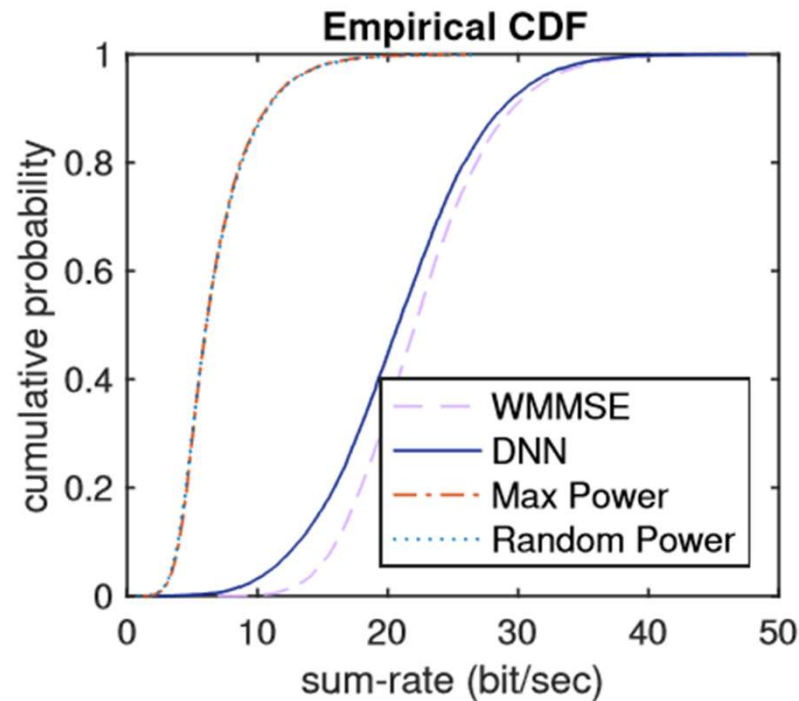
such that the relation below holds true

$$\max_{h \in \mathcal{H}} \max_i |(v(h)_i^T)^2 - NET(h, v^0)_i| \leq \epsilon \quad (15)$$

- If 10 users, $T=20$ and $K=10$, we need **400** layers of neural network and each layer contains **40000** neurons.

In practice, unreasonably effective

- A 3-layer neural network with [200,80,80] neurons in each layer achieves 97% sum rate compared to training labels



The problem is solved



Limitation – Poor Scalability

- Performance deteriorates dramatically when the network size becomes large.

# of users (K)	average sum-rate (bit/sec.)		
	DNN	WMMSE	DNN/WMMSE
10	2.770	2.817	98.33%
20	3.363	3.654	92.04%
30	3.498	4.150	84.29%

Other Limitations of “End-to-end” Learning

- **Huge amounts of samples**
 - 20,000 ~ 100,000,000 samples
 - Optimal labels are difficult to generate
- **Cannot outperform labels**
 - Performance limited by the (sub-optimal) algorithm to generate samples
- **Difficulty in constrained problems**
 - Neural networks are unaware of constraints
- **Weak generalization**
 - Output dimension of neural networks must be fixed

Another approach: Unsupervised Learning

- Consider sum rate maximization

$$\begin{aligned} & \underset{\mathbf{p}}{\text{maximize}} && \sum_{k=1}^K w_k \log_2 \left(1 + \frac{|h_{kk}|^2 p_k}{\sum_{i \neq k} |h_{ki}|^2 p_i + \sigma_k^2} \right) \\ & \text{subject to} && 0 \leq p_k \leq P_{\max}, \forall k, \end{aligned}$$

- Use a neural network to learn the mapping from channel matrix to power $\mathbf{p} = f_{\text{NN}}(\mathbf{H}, \mathbf{w})$
- The **empirical loss function** is

$$\underset{\mathbf{p}}{\text{minimize}} \quad - \sum_{j=1}^N \sum_{k=1}^K w_k \log_2 \left(1 + \frac{|h_{kk}|^2 f_{\text{NN}}(\mathbf{H}_j, \mathbf{w}_j)_k}{\sum_{i \neq k} |h_{ki}|^2 f_{\text{NN}}(\mathbf{H}_j, \mathbf{w}_j)_i + \sigma_k^2} \right)$$

- N is the number of samples

Supervised vs. Unsupervised

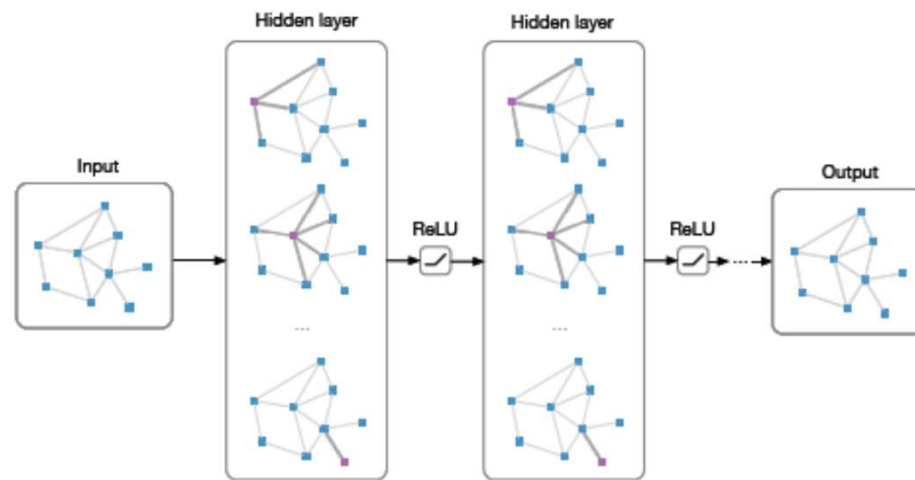
- Supervised
 - Need a traditional algorithm to generate labels
 - MSE loss is often used as loss function
 - Easy to train
- Unsupervised
 - Only channel data is needed (without labels)
 - The objective function of optimization is used as the function of neural network
 - **Hard to train**

This Talk

- To improve from “end-to-end” learning
 - Which model to use? MLP? CNN?
 - What to learn?
- 1. Learn to optimize with graph neural networks
 - To exploit network structure
- 2. Learn optimization policy of a specific algorithm
 - To exploit algorithm structure

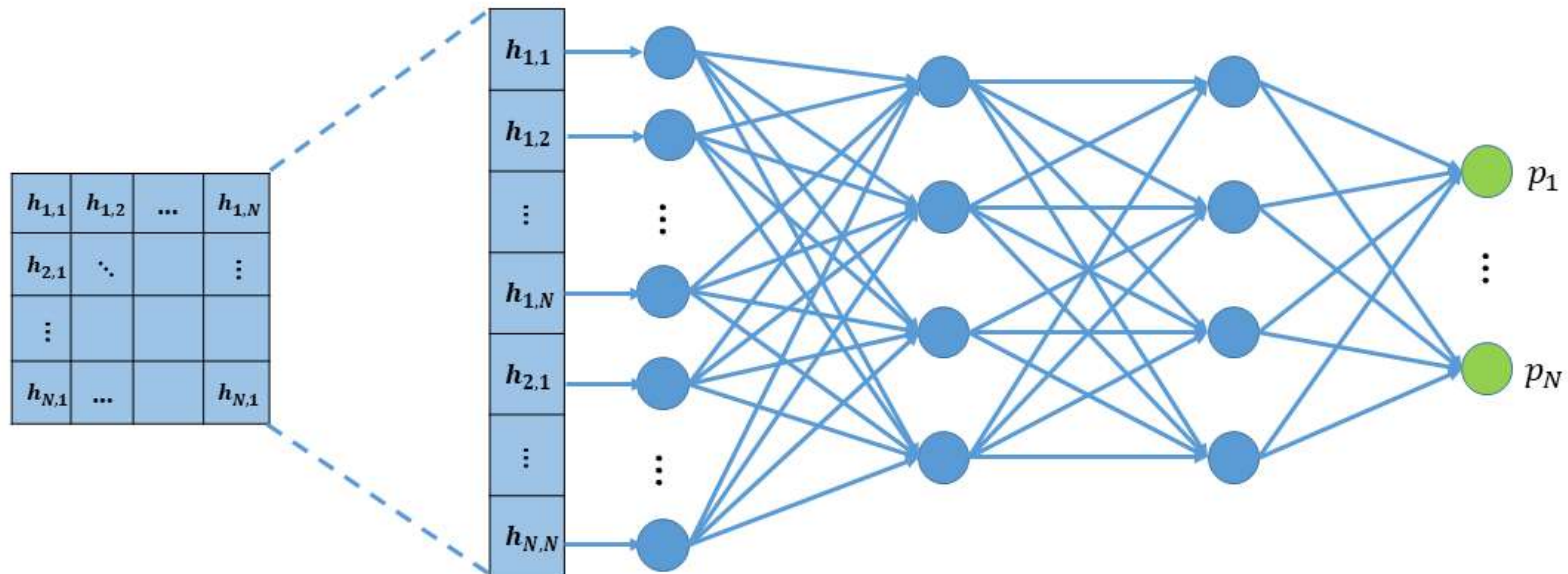
Learn to optimize with graph neural networks

-- exploit network topology



Why MLP is not good enough?

- We flatten the input into vectors before we feed it into the neural network, **thus structure information is lost.**



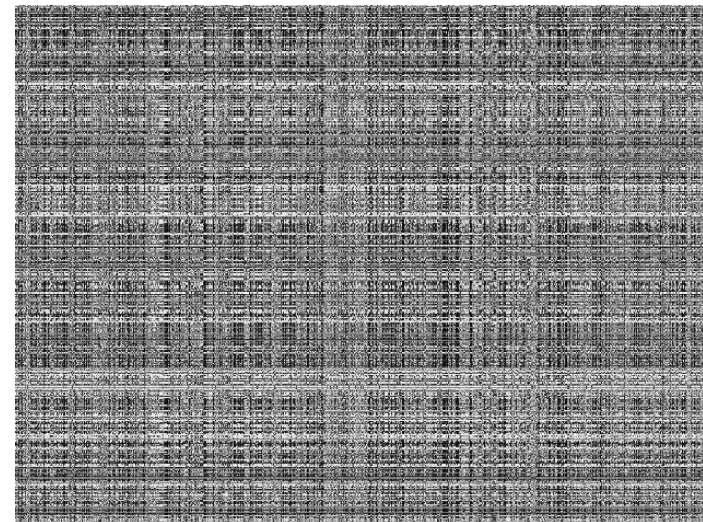
Why MLP is not good enough?

- For example, “structure information” in image processing means the nearby pixels in an image are meaningful
 - But the following two inputs are equivalent for MLP
 - Most of the efforts are spent on discovering the structure

Cat

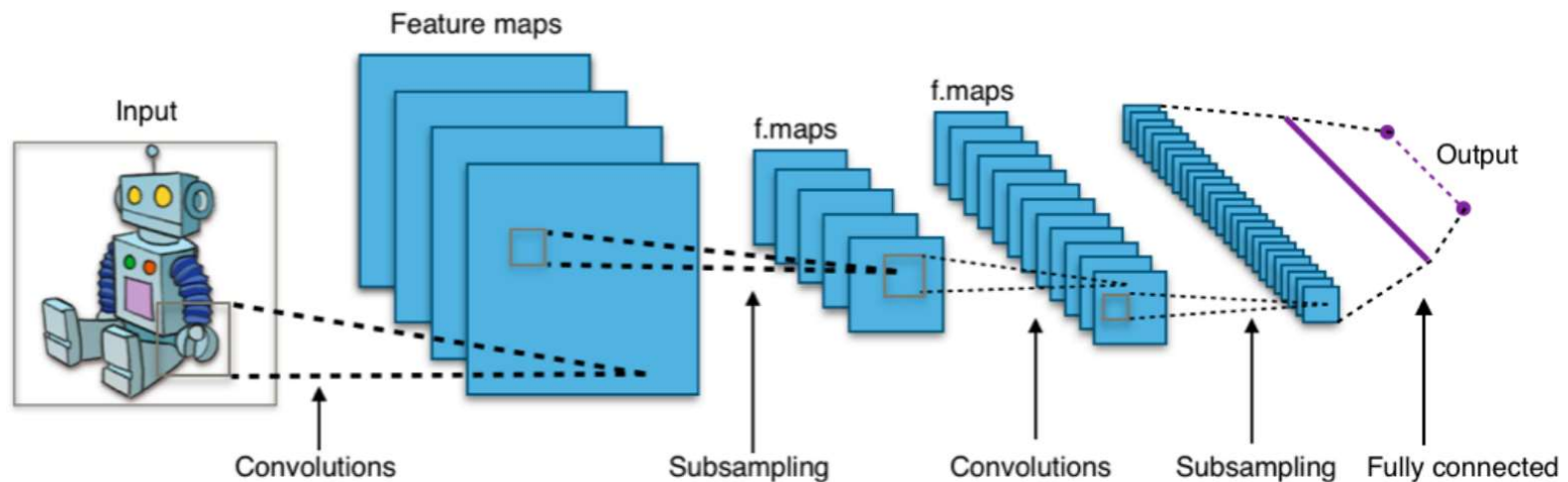


Cat with shuffled pixels



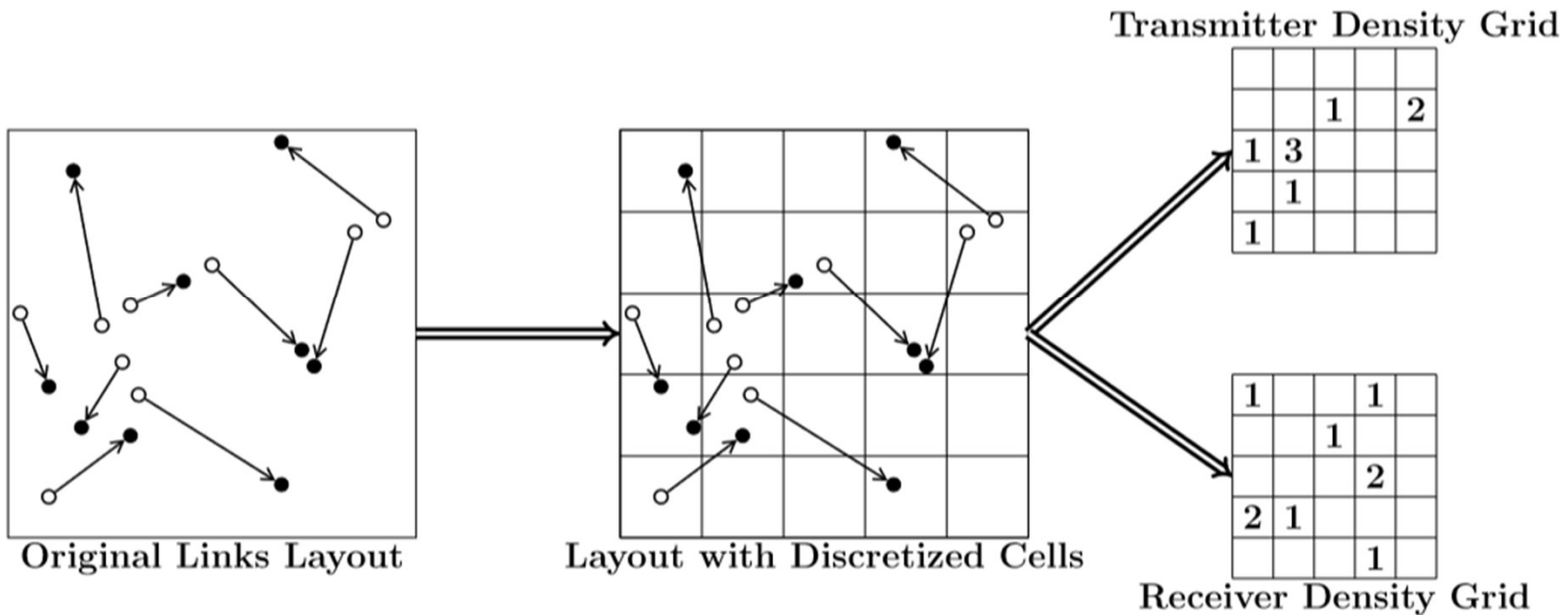
CNN for image processing

- To capture neighborhood information, CNN is proposed for image processing and achieves significant success.
 - It is able to exploit the shift-invariance, local connectivity, and compositionality of image data.



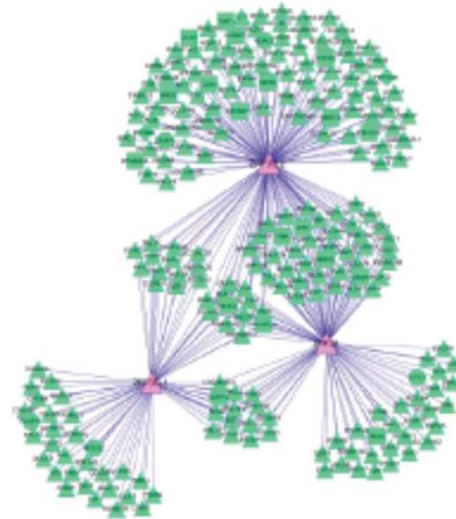
Spatial Convolution for Wireless Communication

- Spatial convolution [Cui19]SAC leverages the geometry of users' locations: nearby users cause the strongest interference.



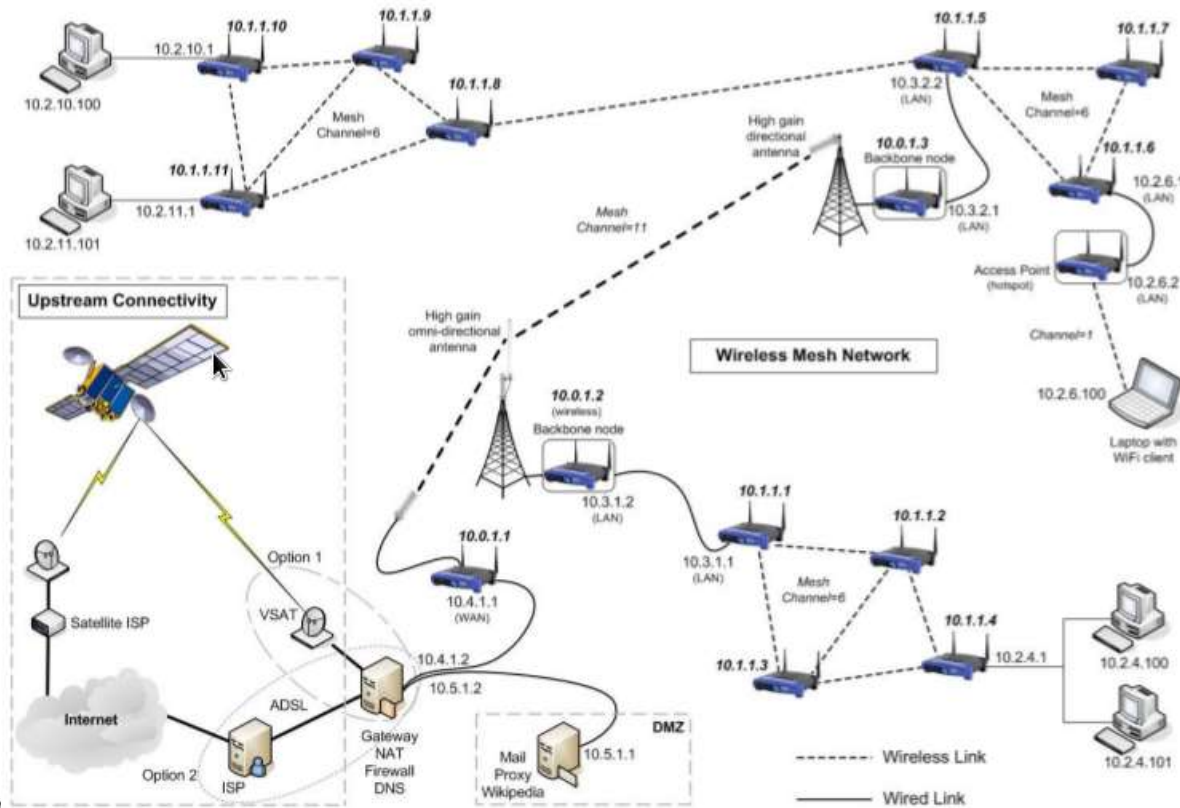
Drawbacks of Spatial Convolution

- CNN is designed for Euclidean data
 - Only **geolocation** can be used as the input, not CSI, distance, or large-scale fading, leading to bad performance when fading exists.
 - Can not directly be applied for **weighted sum rate maximization**.
 - Can not utilize the **topology of the links**.



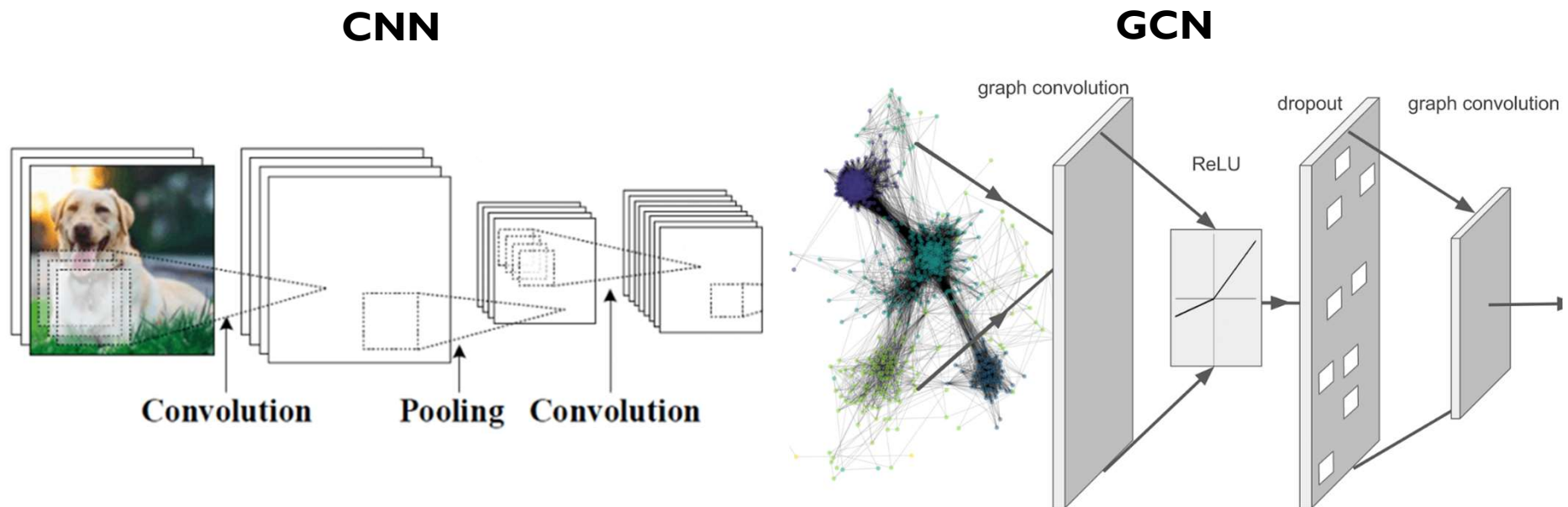
Structure Information for Networks

- Network topology
 - The topology can be naturally modelled as a graph
 - The network topology can be exploited if we use **neural networks on graph**

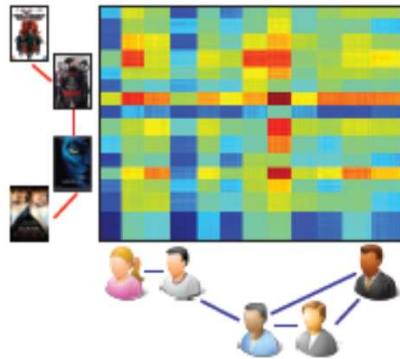


Graph Neural Networks

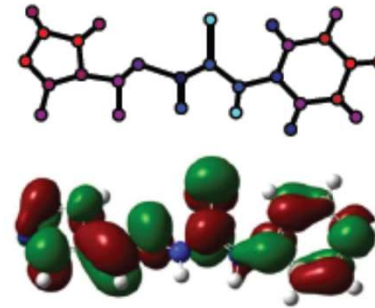
- Like MLP or CNNs, GNNs have layer-wise structures
 - For each layer in CNN, a 2D convolution is applied
 - For each layer in GCN, a graph convolution is applied



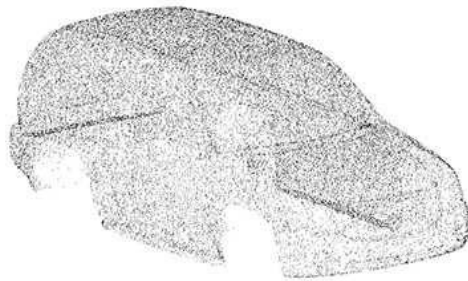
Applications of Graph Neural Networks



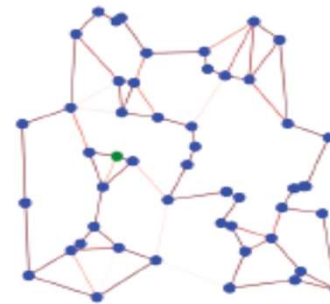
Recommendation



Chemistry



Point clouds



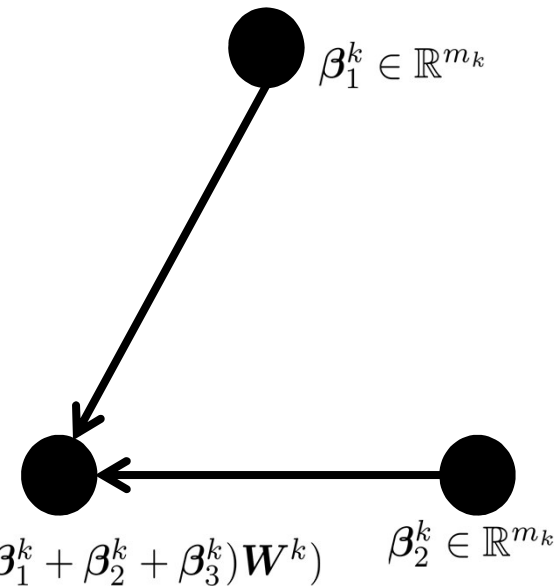
Graph problems

Key Design Steps of GNNs

- In each layer, each node aggregates the information from its neighbors
 - $\beta_v^k \in \mathbb{R}^{m_k}$ denotes the feature vector for node v in k -th NN layer
 - $\mathcal{N}(u)$ denotes the set containing all neighbors of u
 - An **aggregation function** to aggregate features from neighbors
 - A **combination function** to combine the feature of neighbors

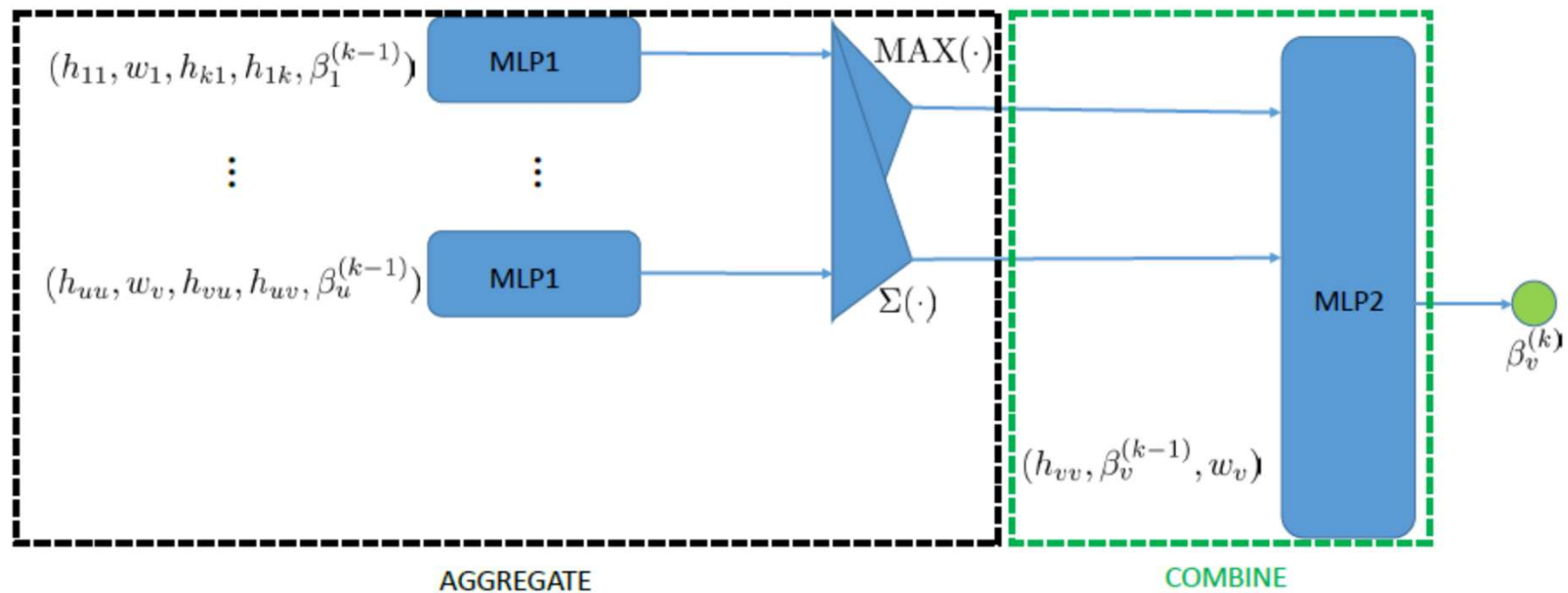
$$\alpha_v^{k+1} = \text{AGGREGATE}^k(\{\beta_u^k : u \in \mathcal{N}(v)\})$$

$$\beta_v^{k+1} = \text{COMBINE}^k(\beta_v^k, \alpha_v^{k+1})$$



Proposed IGCNet

- Interference Graph Convolutional Networks (IGCNet)



- Code available: <https://github.com/yshenaw/Globecom2019>

Motivation of IGCCNet

- Compared with existing works in wireless communications
 - It utilizes network topology
 - It utilizes various kinds of info, e.g., CSI, large-scale fading, distance
 - It achieves stable performance with varying network sizes
- Compared with existing works of GNNs
 - Existing works of GNN can not deal with edge features, e.g., GCN [Kipf17ICLR], GIN [Xu19ICLR], structure2Vec [Dai16ICML],
 - or too complicated and slow, e.g., NRI [Kipf18ICML], GAT [Shang18arxiv]

Simulations

- Different methods
 - **IGCNet**: Proposed method
 - **WMMSE**: The most popular optimization-based method [Shi11TSP]
 - **MLP**: use MLP to approximate WMMSE [Sun18TSP]
 - **DPC**: use CNN to approximate WMMSE [Lee18CL]
 - **PCNet**: MLP with unsupervised training [Liang18arxiv]
 - **Baseline**: Activate pairs with largest channel gains, the simplest method and ignoring the interference
 - Note: we omit methods that can only handle geolocation or distance inputs [Cui19SAC][Lee19arxiv]

Simulations – Scalability

- Maintaining the performance when the network size grows
 - The decision is made locally at each node, less impact from the scale
 - K : number of users

TABLE II
AVERAGE SUM RATE UNDER EACH SETTING. THE RESULTS ARE
NORMALIZED BY THE SUM RATE ACHIEVED BY WMMSE.

	IGCNet	MLP	PCNet	DPC	Baseline
$K = 10$	102.6%	98.2%	101.4%	95.1%	89.1%
$K = 20$	102.7%	92.3%	90.2%	83.1%	86.6%
$K = 30$	102.4%	85.3%	87.6%	79.3%	84.4%

Simulations – Robustness

- Robustness to imperfect CSI

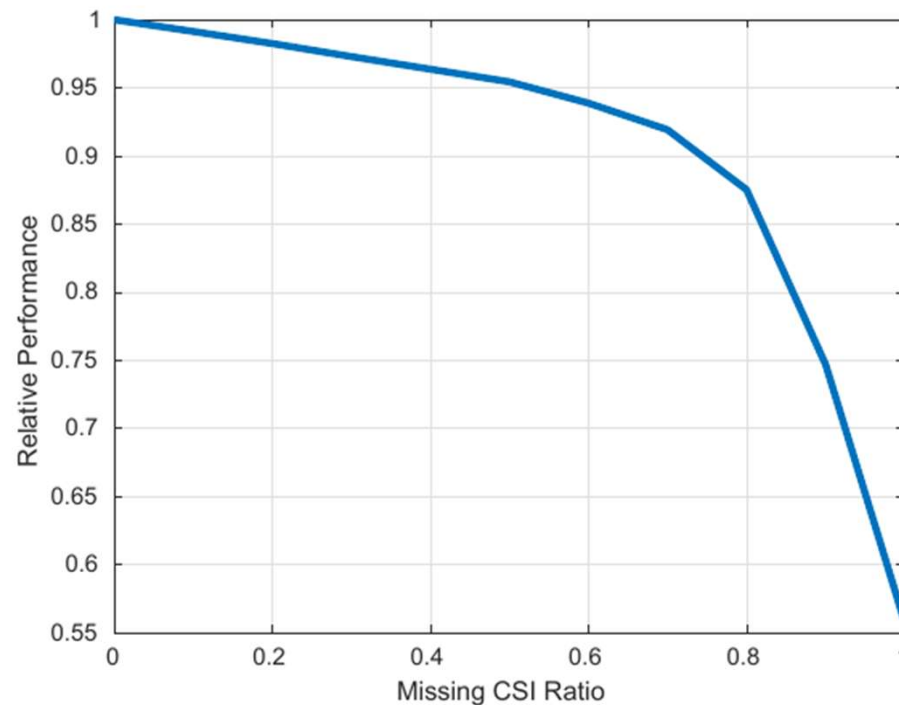


Fig. 5. The relative performance versus the missing ratio of CSI.

Simulations – Computation Efficiency

- Fast computation
 - WMMSE involves many iterations, and each iteration is $O(K^2)$
 - The total complexity of neural networks is $O(K^2)$

TABLE IV
AVERAGE RUNNING TIME FOR THE ALGORITHMS UNDER EACH SETTING
(IN MILLISECONDS).

	$K = 10$	$K = 20$	$K = 30$
IGCNet	0.14ms	0.27ms	0.48ms
WMMSE	9.31ms	24.1ms	31.4ms

Benefits of IGCCNet

- Maintaining the performance when the network size grows
 - Scalable
- Able to handle non-Euclidean features
 - Incorporate CSI and handle weighted objective functions
- Robust to missing data
 - Robust to CSI uncertainty
- Fast computation
 - Real-time execution
- Few training samples without labels
 - Easy to implement

LORM: Learn to optimize with few samples



Learning for Mixed Integer Nonlinear Program

- Mixed Integer Nonlinear Programming (MINLP) Problems

$$\begin{array}{ll} \underset{\mathbf{w}, \mathbf{a}}{\text{minimize}} & f(\mathbf{a}, \mathbf{w}) \\ \text{subject to} & Q(\mathbf{a}, \mathbf{w}) \leq 0 \\ & a_i \in \mathbb{N}, w_i \in \mathbb{C}. \end{array}$$

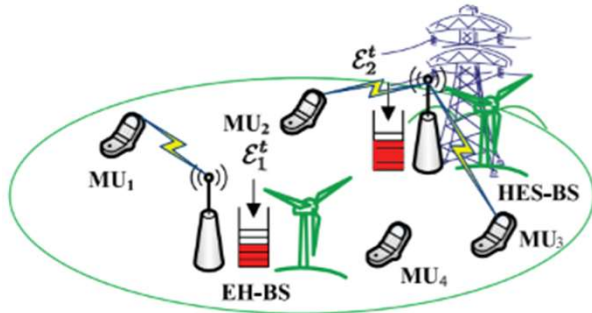
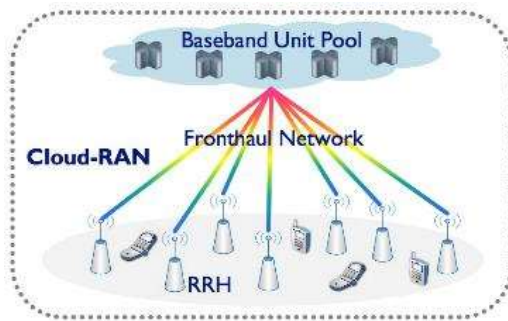
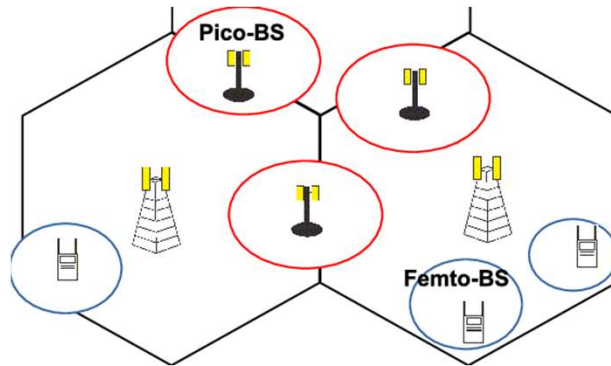
Diagram illustrating the components of a Mixed Integer Nonlinear Program (MINLP) problem:

- $f(\mathbf{a}, \mathbf{w})$ is labeled as the **Cost function**.
- $Q(\mathbf{a}, \mathbf{w}) \leq 0$ is labeled as **Resource constraints**.
- $a_i \in \mathbb{N}$ is labeled as **Integer variables**.
- $w_i \in \mathbb{C}$ is labeled as **Continuous variables**.

- NP-hard because of the combinatorial variables
- Both continuous and discrete variables

- Many resource allocation problems are MINLP

Typical MINLP Resource Management Problems



- User Association in HetNets
 - Q. Ye, B. Rong, Y. Chen, M. Al-Shalash, C. Caramanis, and J. G. Andrews, “User association for load balancing in heterogeneous cellular networks,” *IEEE Trans. Wireless Commun.*, vol. 12, pp. 2706–2716, Jun. 2013.

- Power Minimization in Cloud RANs
 - Y. Shi, J. Zhang, and K. B. Letaief, “Group sparse beamforming for green Cloud-RAN,” *IEEE Trans. Wireless Commun.*, vol. 13, no. 5, pp. 2809-2823, May 2014.

- Computation off-loading in MEC
 - Y. Mao, J. Zhang, and K. B. Letaief, “Dynamic computation offloading for mobile-edge computing with energy harvesting devices,” *IEEE J. Sel. Areas Commun.*, vol. 34, pp. 3590–3605, Dec. 2016.

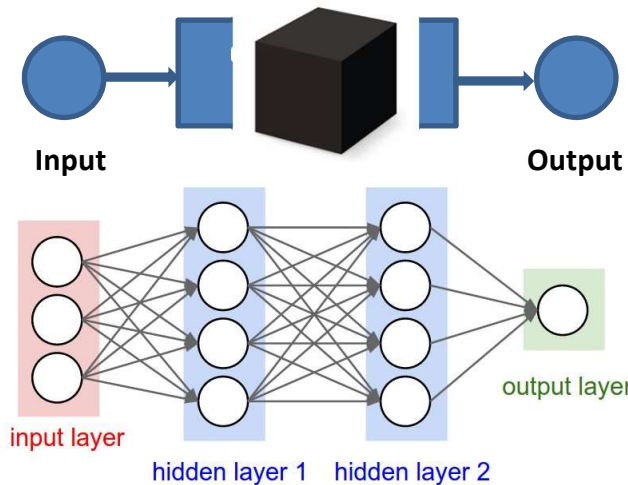
Basic Algorithmic Approaches

- **Global optimization algorithms**
 - exponential time complexity
 - only work for very small problems
- **Heuristic algorithms**
 - Examples: **greedy algorithm** for user selection or **sub-optimal algorithm** like zero-forcing
 - hard to design good ones
 - non-negligible gap to the optimal solution
 - difficult to meet real-time requirement

Different ways to “learn to optimize”

- **End-to-end-learning**

- Directly learn the input-output mapping



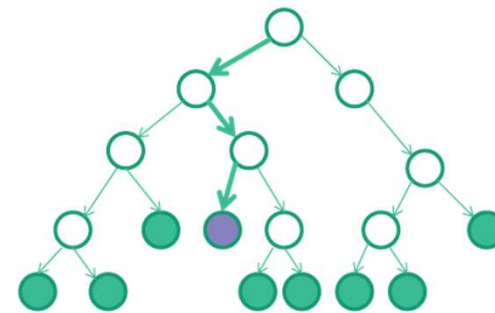
- Needs large model, many samples
- Different to handle constraints

- **Optimization policy learning**

- Learn the optimal policy in a specific algorithm



Branch-and-Bound

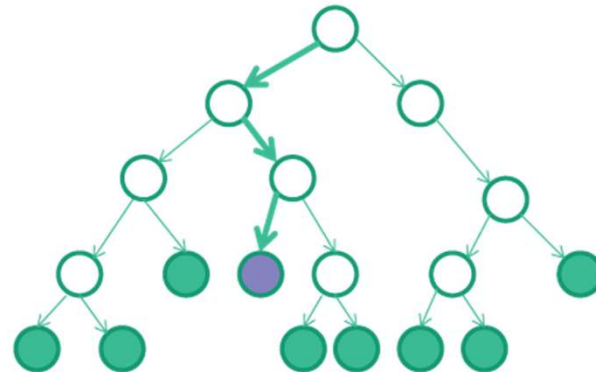


- Exploits algorithm structure, requires few samples
- Capable to handle constraints

Branch-and-Bound

- A global optimization algorithm for MINLP

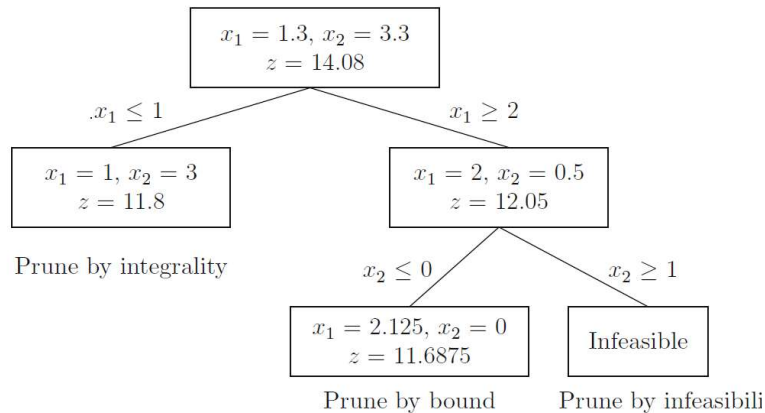
Branch-and-Bound



- Three policies:
 - **Node selection policy**: select a node in binary search tree
 - **Variable selection policy**: select a variable to branch
 - **Prune policy**: whether to expand two children or not (time is saved if not expand)

Branch-and-Bound

$$\begin{aligned}
 \max \quad & 5.5x_1 + 2.1x_2 \\
 & -x_1 + x_2 \leq 2 \\
 & 8x_1 + 2x_2 \leq 17 \\
 & x_1, x_2 \geq 0 \\
 & x_1, x_2 \text{ integer.}
 \end{aligned}$$



Branch-and-Bound Algorithm

0. Initialize

$$\mathcal{L} := \{N_0\}, \underline{z} := -\infty, (x^*, y^*) := \emptyset.$$

1. Terminate?

If $\mathcal{L} = \emptyset$, the solution (x^*, y^*) is optimal.

2. Select node

Choose a node N_i in \mathcal{L} and delete it from \mathcal{L} .

3. Bound

Solve LP_i . If it is infeasible, go to Step 1. Else, let (x^i, y^i) be an optimal solution of LP_i and z_i its objective value.

4. Prune

If $z_i \leq \underline{z}$, go to Step 1.

If (x^i, y^i) is feasible to MILP, set $\underline{z} := z_i$, $(x^*, y^*) := (x^i, y^i)$ and go to Step 1.

Otherwise:

5. Branch

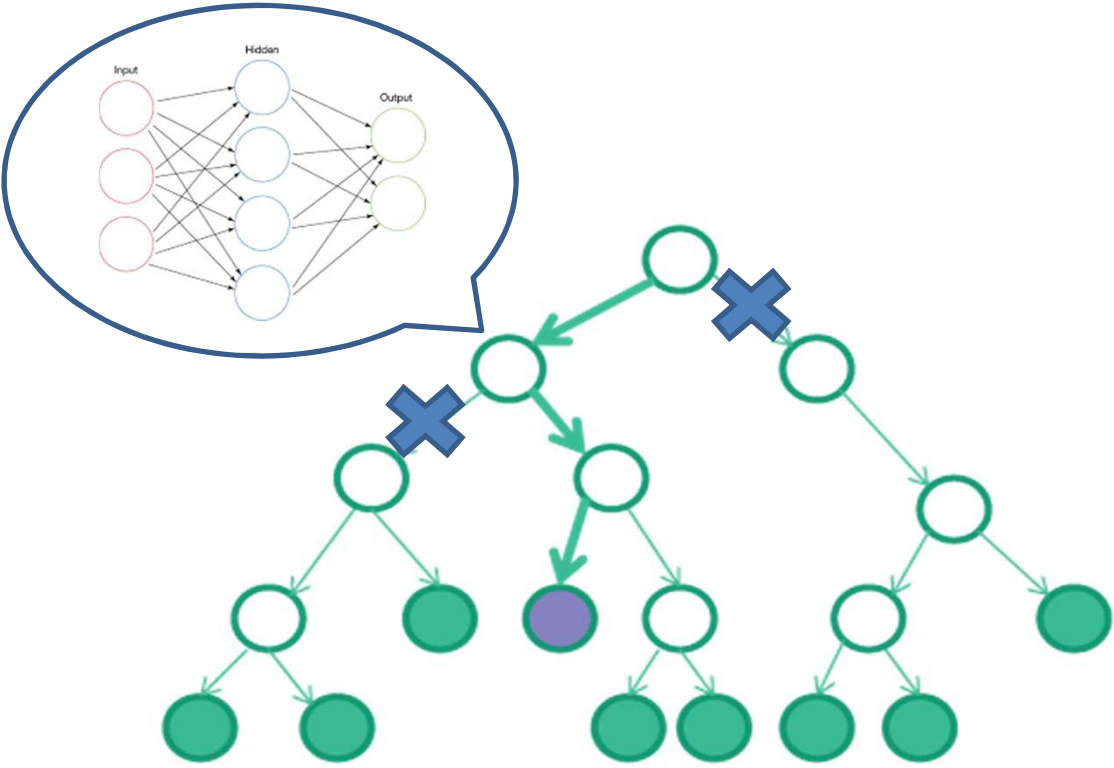
From LP_i , construct $k \geq 2$ linear programs $LP_{i_1}, \dots, LP_{i_k}$ with smaller feasible regions whose union does not contain (x^i, y^i) , but contains all the solutions of LP_i with $x \in \mathbb{Z}^n$. Add the corresponding new nodes N_{i_1}, \dots, N_{i_k} to \mathcal{L} and go to Step 1.

[Conforti, et al., 2014]

Learn to Prune

- Prune policies
 - Prune by bound: the lower bound is worse than best integer solution
 - Prune by infeasibility: the relaxed problem is infeasible
 - Prune by integrality: the relaxed problem has an integer solution
- Most of the time spent on checking non-optimal nodes
- **Insight:** We want a good enough solution rather than optimality guarantee
 - **Learning Pruning Policy** - Classification

Learn to Prune



Imitation Learning

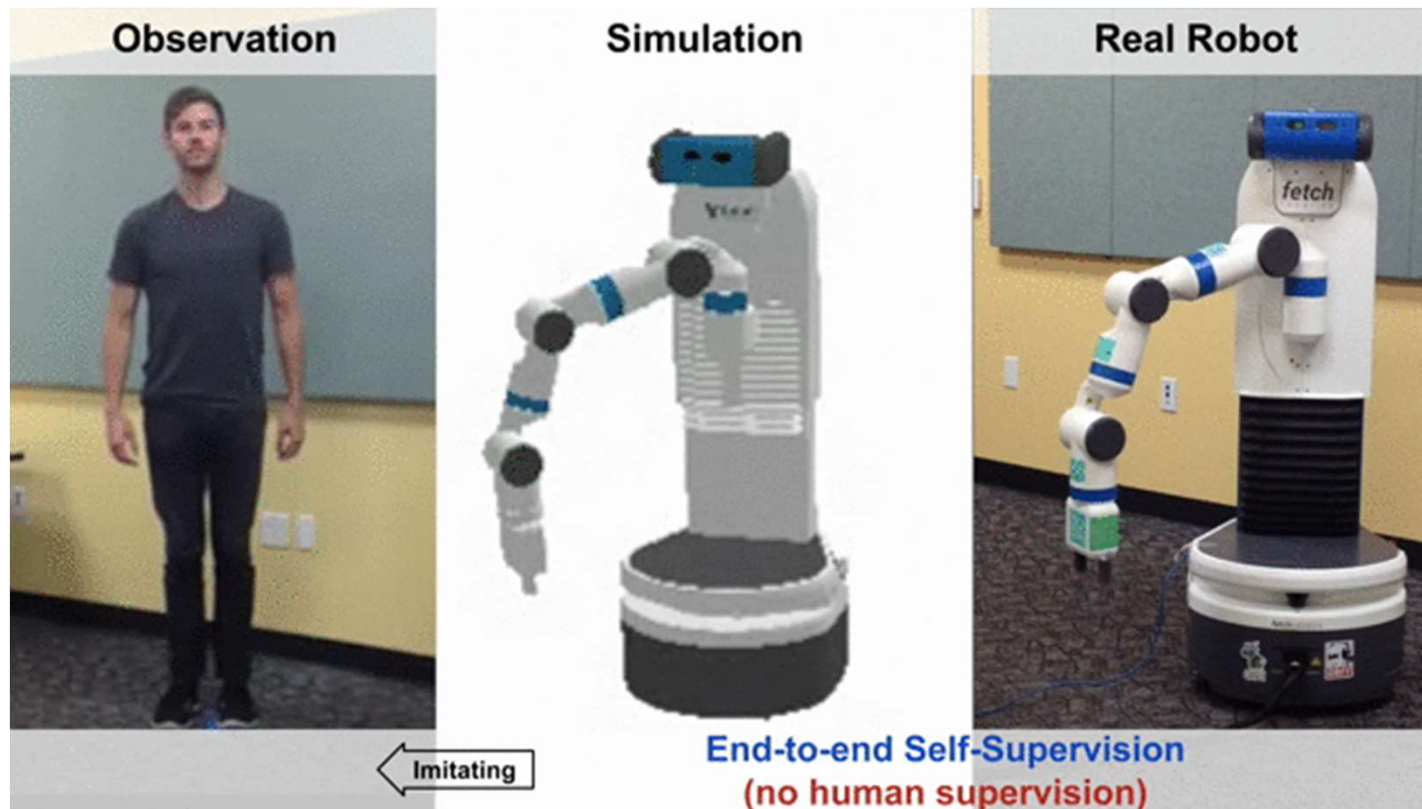
- Imitation learning to learn optimization policy
 - Imitation learning mimics an optimal policy
 - In learning the pruning policy, the optimal policy only preserves nodes containing the optimal solution
- Vs supervised learning
 - Iteratively collect new data and training – **better generalization**
- Vs reinforcement learning
 - Learn directly from the optimal policy rather than indirectly from the reward – **lower sample complexity**

Foundations and Trends® in Robotics
Vol. 7, No. 1-2 (2018) 1–179
© 2018 T. Osa, J. Pajarinen, G. Neumann,
J. A. Bagnell, P. Abbeel and J. Peters
DOI: 10.1561/23000000053

An Algorithmic Perspective on Imitation Learning

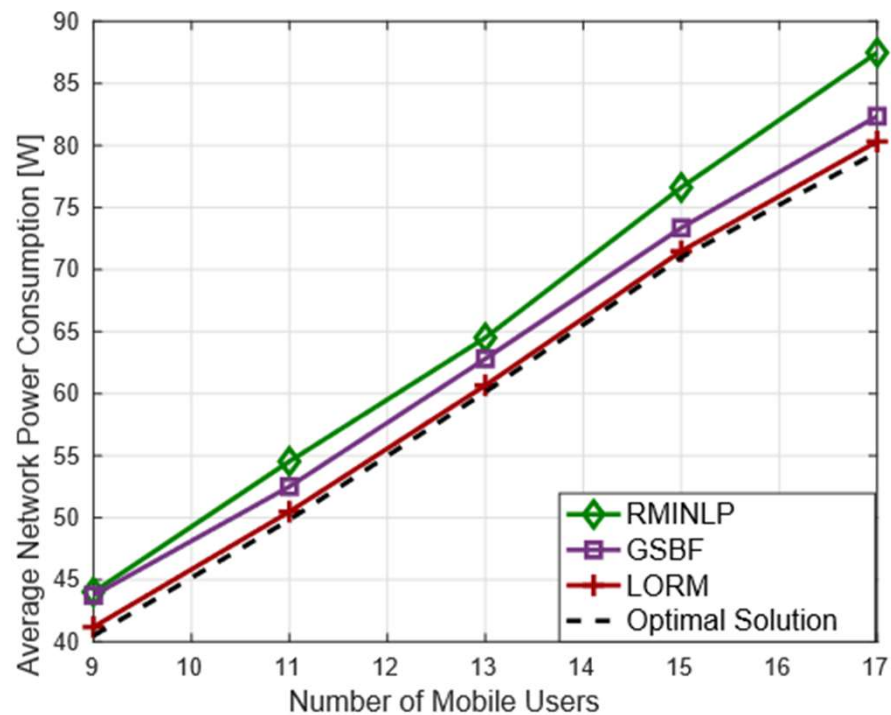
Imitation Learning

- Mimic an optimal behavior



Simulation Results – Better Performance

Network power minimization in C-RANs



- **GSBF**: State-of-art
- **RMINLP**: Heuristic method
- **Branch-and-Bound**
- **LORM**: Proposed method
 - 50 training samples
 - Near optimal solution

Simulation Results – Computation Speedup

Setting	Branch-and-Bound	LORM	GSBF	RMINLP
$L = 10, K = 13, \text{TSINR} = 4$	91.76s	0.892s	2.562s	5.264s
$L = 10, K = 15, \text{TSINR} = 4$	96.40s	1.157s	3.680s	8.136s
$L = 10, K = 17, \text{TSINR} = 4$	142.0s	2.920s	5.474s	12.64s

- Use only 50 problem instances – easy for training
- 70x speedup to the branch-and-bound.
- 2x speedup to the state-of-the-art method.

Simulation Results – Generalization

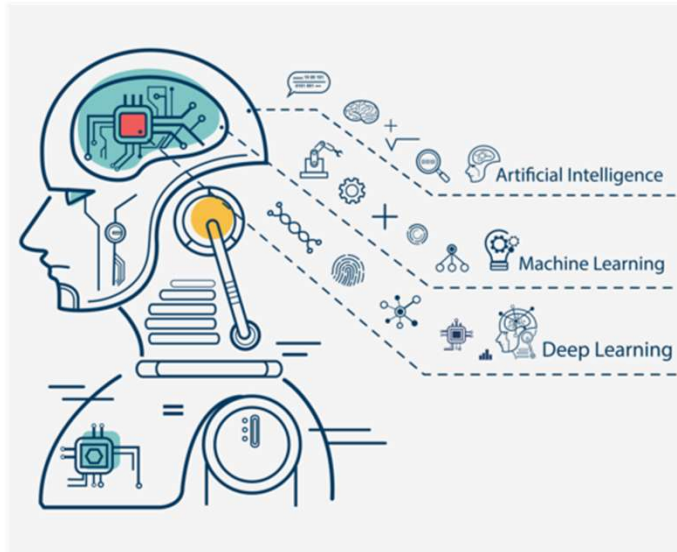
The Gap to the Optimal Objective Value

Setting	LORM (full training)	LORM (train on $L = 6, K = 9, \text{TSINR} = 0$)	RMINLP	GSBF
$L = 10, K = 7, \text{TSINR} = 4$	1.87%	4.39%	7.94%	12.9%
$L = 10, K = 9, \text{TSINR} = 4$	1.73%	2.97%	8.71%	8.04%
$L = 10, K = 11, \text{TSINR} = 4$	1.30%	1.72%	9.44%	5.36%
$L = 10, K = 13, \text{TSINR} = 4$	1.41%	1.41%	7.27%	4.45%
$L = 10, K = 15, \text{TSINR} = 4$	0.70%	1.48%	7.94%	3.36%

Advantages of LORM

- Near-optimal performance with **few training samples**
 - By learning the optimization policy via imitation learning
- It **outperforms** the non-optimal labels
 - Obtaining a better solution by pruning fewer nodes
- It guarantees **feasibility** of constraints
 - Retain algorithm structure
- It **generalizes** to different system configurations, and is able to scale up to larger problem sizes
 - The input and output dimensions of the pruning policy are invariant to problem sizes

Conclusions



Conclusions

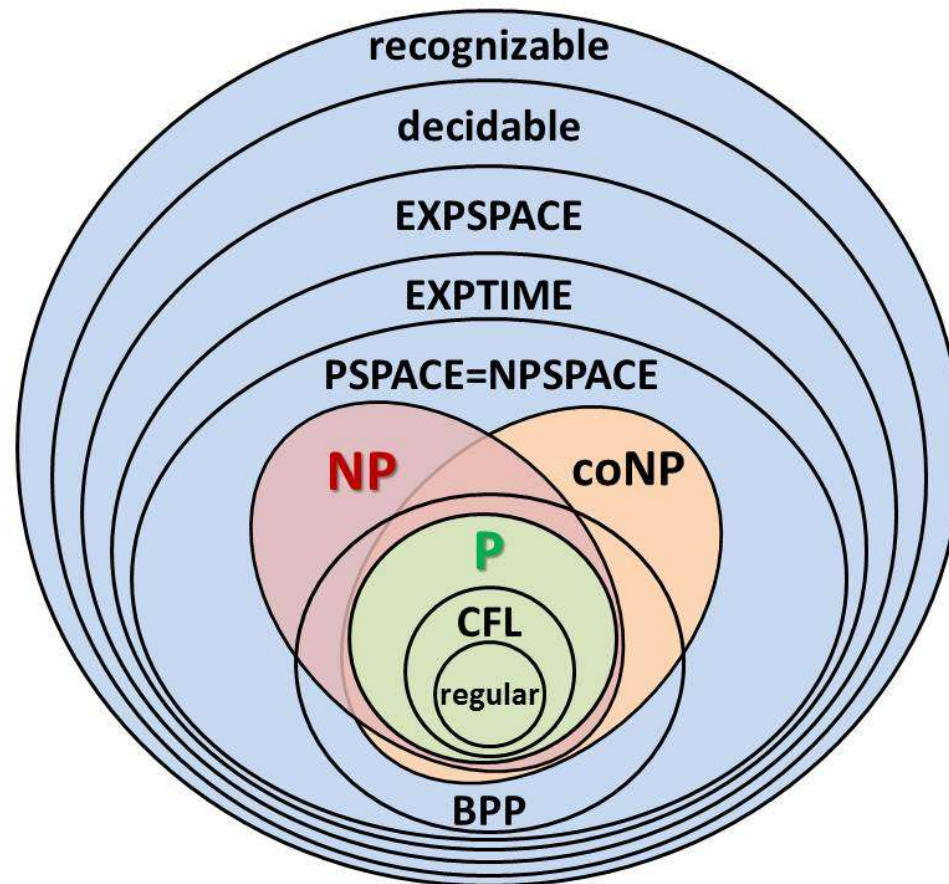
- Learn to optimize with network topology
 - Graph neural network (GNN) based approach
 - Scalable to large network sizes
 - Unsupervised (no need for labelling)
- Learn optimization policies
 - Able to handle general MINLP problems
 - Supervised (few labelled samples)
 - Able to outperform labels
 - Constraint guarantee
 - Good generalization

Conclusions

- Great potentials of “learn to optimize” for wireless
 - Higher computational efficiency
 - Avoid hand-crafted algorithm design
 - Close-to-optimal performance
- Key takeaways
 - End-to-end learning is not efficient
 - Do not directly apply MLP or CNN
 - Exploit structures
 - Network structure (IGCNet)
 - Algorithm structure (LORM)

One Picture to Recall

- A big universe of problems



One Slide to Take Away

Learn to Optimize (L2O)

Learn with
Network Topology

Learn
Optimization
Policy

1. Figure out what problems can be solved via L2O
2. Find effective “simple” models to solve them
 - Few samples
 - Easy to train
 - High computational efficiency

References on “learn to optimize”

- B. Yoshua, L. Andrea, and A. Prouvost, “Machine learning for combinatorial optimization: a methodological tourd’horizon,” arXiv preprint arXiv:1811.06128, 2018.
- M.-F. Balcan, T. Dick, T. Sandholm, and E. Vitercik, “Learning to branch,” in *Proc. Int. Conf. Mach. Learning*, vol. 80, pp. 344–353, Jul. 2018.
- J. Song, R. Lanka, A. Zhao, Y. Yue, and M. Ono, “Learning to search via self-imitation with application to risk-aware planning,” in *Proc. Adv. Neural Inform. Process. Syst. Workshop*, Dec. 2017.
- H. He, H. Daume III, and J. M. Eisner, “Learning to search in branch and bound algorithms,” in *Proc. Adv. Neural Inform. Process. Syst.*, pp. 3293–3301, Dec. 2014.

References (ML in Wireless)

- H. Sun, X. Chen, Q. Shi, M. Hong, X. Fu, and N. D. Sidiropoulos, “Learning to optimize: Training deep neural networks for interference management,” *IEEE Trans. Signal Process.*, vol. 66, pp. 5438 – 5453, Oct. 2018.
- W. Lee, M. Kim, and D.-H. Cho, “Deep power control: Transmit power control scheme based on convolutional neural network,” *IEEE Commun. Lett.*, vol. 22, pp. 1276–1279, Apr. 2018.
- F. Liang, C. Shen, W. Yu, and F. Wu, “Towards optimal power control via ensembling deep neural networks,” arXiv preprint arXiv:1807.10025, 2018.
- W. Cui, K. Shen, and W. Yu, “Spatial Deep Learning for Wireless Scheduling”, *IEEE Journal on Selected Areas in Communications*, 2019.
- H. Lee, S. H. Lee, T. Q. S. Quek, “Deep Learning for Distributed Optimization: Applications to Wireless Resource Management”, *IEEE Journal on Selected Areas in Communications*, 2019.
- M. Lee, G. Yu, and G. Y. Li, “Graph embedding based wireless link scheduling with few training samples,” arXiv preprint arXiv:1906.02871, 2019.

References on Graph Neural Networks

- Survey

- Z. Wu, S. Pan, F. Chen, G. Long, C. Zhang, and P. S. Yu, “A comprehensive survey on graph neural networks,” arXiv preprint arXiv:1901.00596, 2019.
- K. Xu, W. Hu, J. Leskovec, and S. Jegelka, “How powerful are graph neural networks?,” Proc. Int. Conf. Learning Representation, May 2019.

- GNNs

- T. N. Kipf and M. Welling, “Semi-supervised classification with graph convolutional networks,” Proc. Int. Conf. Learning Representation, Apr. 2017.
- P. Veličković, G. Cucurull, A. Casanova, etc “Graph Attention Networks,” Proc. Int. Conf. Learning Representation, Apr. 2018.

- Applications

- R. Ying, R. He, K. Chen, P. Eksombatchai, W. L. Hamilton, and J. Leskovec, “Graph convolutional neural networks for web-scale recommender systems,” in Proc. ACM Int. Conf. Knowl. Discovery Data Mining, pp. 974–983, ACM, Aug. 2018.
- Z. Li, Q. Chen, and V. Koltun, “Combinatorial optimization with graph convolutional networks and guided tree search,” in Proc. Adv. Neural Inform. Process. Syst., pp. 539–548, Dec. 2018.
- G. Li, M. Müller, A. Thabet, and B. Ghanem, “Can GCNs Go as Deep as CNNs?” in Proc. Int. Conf. Comput. Vision, Oct. 2019.

References (Our Results)

- Y. Shen, Y. Shi, **J. Zhang**, and K. B. Letaief, “LORA: Learning to optimize for resource allocation in wireless networks with few training samples,” submitted.
- Y. Shen, Y. Shi, **J. Zhang**, and K. B. Letaief, “Transfer learning for mixed-integer resource allocation problems in wireless networks,” in *Proc. IEEE Int. Conf. Commun. (ICC)*, Shanghai, China, May 2019.
- Y. Shen, Y. Shi, **J. Zhang**, and K. B. Letaief, “Scalable network adaption for green Cloud-RANs: An imitation learning approach,” in *Proc. IEEE Global Conf. Signal and Inf. Process. (GlobalSIP)*, Anaheim, CA, Nov. 2018.
- Y. Shen, Y. Shi, **J. Zhang**, and K. B. Letaief, “A Graph Neural Network Approach for Scalable Wireless Power Control,” IEEE GLOBECOM Workshop, Dec. 2019.

Thank you!

- For more details

<http://www.eie.polyu.edu.hk/~jeiezhang/>