

Machine learning for discrete optimization: Theoretical guarantees and applied frontiers

Ellen Vitercik
Stanford University

How to integrate **machine learning** into **discrete optimization**?

- **Algorithm configuration**
How to tune an algorithm's parameters?
- **Algorithm selection**
Given a variety of algorithms, which to use?
- **Algorithm design**
Can machine learning guide algorithm discovery?

How to integrate **machine learning** into **discrete optimization**?



Algorithm configuration

How to tune an algorithm's parameters?



Algorithm selection

Given a variety of algorithms, which to use?



Algorithm design

Can machine learning guide algorithm discovery?

Algorithm configuration

Example: **Integer programming solvers**

Most popular tool for solving combinatorial (& nonconvex) problems



Routing



Manufacturing



Scheduling



Planning



Finance

Algorithm configuration

IP solvers (CPLEX, Gurobi) have a **ton** parameters

- CPLEX has **170-page** manual describing **172** parameters
- Tuning by hand is notoriously **slow, tedious,** and **error-prone**

CPX_PARAM_NODEFILEIND 100	CPX_PARAM_TRELIM 160	CPX_PARAM_RANDOMSEED 130	CPXPARAM_MIP_Pool_RelGap 148	CPX_PARAM_FLOWCOVERS 70	CPX_PARAM_BRDIR 39
CPX_PARAM_NODELIM 101	CPX_PARAM_TUNINGDETTILIM 160	CPX_PARAM_REDUCE 131	CPXPARAM_MIP_Pool_Replace 151	CPX_PARAM_FLOWPATHS 71	CPX_PARAM_BTTL 40
CPX_PARAM_NODESEL 102	CPX_PARAM_TUNINGDISPLAY 162	CPX_PARAM_REINV 131	CPXPARAM_MIP_Strategy_Branch 39	CPX_PARAM_FPHEUR 72	CPX_PARAM_CALCQCPCDUALS 41
CPX_PARAM_NUMERICALEMPHASIS 102	CPX_PARAM_TUNINGMEASURE 163	CPX_PARAM_RELAXPREIND 132	CPXPARAM_MIP_Strategy_MIQCPStrat 93	CPX_PARAM_FRACCAND 73	CPX_PARAM_CLIQUES 42
CPX_PARAM_NZREADLIM 103	CPX_PARAM_TUNINGREPEAT 164	CPX_PARAM_RELOBJDIF 133	CPXPARAM_MIP_Strategy_StartAlgorithm 139	CPX_PARAM_FRACCUTS 73	CPX_PARAM_CLOCKTYPE 43
CPX_PARAM_OBJDIF 104	CPX_PARAM_TUNINGTILIM 165	CPX_PARAM_REPAIRTRIES 133	CPXPARAM_MIP_Strategy_VariableSelect 166	CPX_PARAM_FRACPASS 74	CPX_PARAM_CLONELOG 43
CPX_PARAM_OBJLLIM 105	CPX_PARAM_VARSEL 166	CPX_PARAM_REPEATPRESOLVE 134	CPXPARAM_MIP_SubMIP_NodeLimit 155	CPX_PARAM_GUBCOVERS 75	CPX_PARAM_COEREDIND 44
CPX_PARAM_OBJULIM 105	CPX_PARAM_WORKDIR 167	CPX_PARAM_RINSHEUR 135	CPXPARAM_OptimalityTarget 106	CPX_PARAM_HEURFREQ 76	CPX_PARAM_COLREADLIM 45
CPX_PARAM_PARALLELMODE 108	CPX_PARAM_WORKMEM 168	CPX_PARAM_RLT 136	CPXPARAM_Output_WriteLevel 169	CPX_PARAM_IMPLBD 76	CPX_PARAM_CONFLICTDISPLAY 46
CPX_PARAM_PERIND 110	CPX_PARAM_WRITELEVEL 169	CPX_PARAM_ROWREADLIM 141	CPXPARAM_Preprocessing_Aggregator 19	CPX_PARAM_INTSOLFILEPREFIX 78	CPX_PARAM_COVERS 47
CPX_PARAM_PERLIM 111	CPX_PARAM_ZEROHALFCUTS 170	CPX_PARAM_SCAIND 142	CPXPARAM_Preprocessing_Fill 19	CPX_PARAM_INTSOLLIM 79	CPX_PARAM_CPUMASK 48
CPX_PARAM_POLISHAFTERDETTIME 111	CPXPARAM_Benders_Strategy 30	CPX_PARAM_SCRIND 143	CPXPARAM_Preprocessing_Linear 120	CPX_PARAM_ITLIM 80	CPX_PARAM_CRAIN 50
CPX_PARAM_POLISHAFTEREPAGAP 112	CPXPARAM_Benders_Tolerances_feasibilitycut 35	CPX_PARAM_SIFTALG 143	CPXPARAM_Preprocessing_Reduce 131	CPX_PARAM_LANDPCUTS 82	CPX_PARAM_CUTLO 51
CPX_PARAM_POLISHAFTEREPGAP 113	CPXPARAM_Benders_Tolerances_optimalitycut 36	CPX_PARAM_SIFTDISPLAY 144	CPXPARAM_Preprocessing_Symmetry 156	CPX_PARAM_LBHEUR 81	CPX_PARAM_CUTPASS 52
CPX_PARAM_POLISHAFTERINTSOL 114	CPXPARAM_Conflict_Algorithm 46	CPX_PARAM_SIFTTILIM 145	CPXPARAM_Read_DataCheck 54	CPX_PARAM_LPMETHOD 136	CPX_PARAM_CUTSFACTOR 52
CPX_PARAM_POLISHAFTERNODE 115	CPXPARAM_CPUmask 48	CPX_PARAM_SIMDISPLAY 145	CPXPARAM_Read_Scale 142	CPX_PARAM_MFCUTS 82	CPX_PARAM_CUTUP 53
CPX_PARAM_POLISHAFTERTIME 116	CPXPARAM_DistMIP_Rampup_Duration 128	CPX_PARAM_SINGLIM 146	CPXPARAM_ScreenOutput 143	CPX_PARAM_MEMORYEMPHASIS 83	CPXPARAM_DATACHECK 54
CPX_PARAM_POLISHTIME (deprecated) 116	CPXPARAM_LPMethod 136	CPX_PARAM_SOLNPOOLGAP 146	CPXPARAM_Sifting_Algorithm 143	CPX_PARAM_MIPCBREDLP 84	CPX_PARAM_DEPIND 55
CPX_PARAM_POPULATELIM 117	CPXPARAM_MIP_Cuts_BQP 38	CPX_PARAM_SOLNPOOLCAPACITY 147	CPXPARAM_Sifting_Display 144	CPX_PARAM_MIPDISPLAY 85	CPX_PARAM_DETTILIM 56
CPX_PARAM_PPRIND 118	CPXPARAM_MIP_Cuts_LocallyImplied 77	CPX_PARAM_SOLNPOOLREPLACE 151	CPXPARAM_Sifting_Iterations 145	CPX_PARAM_MIPEMPHASIS 87	CPX_PARAM_DISJCUTS 57
CPX_PARAM_PREDUAL 119	CPXPARAM_MIP_Cuts_RLT 136	CPX_PARAM_SOLNPOOLINTENSITY 149	CPXPARAM_Simplex_Display 145	CPX_PARAM_MIPINTERVAL 88	CPX_PARAM_DIVETYPE 58
CPX_PARAM_PREIND 120	CPXPARAM_MIP_Cuts_ZeroHalfCut 170	CPX_PARAM_SOLNPOOLREPLACE 151	CPXPARAM_Simplex_Limits_Singularity 146	CPX_PARAM_MIPKAPPASTATS 89	CPX_PARAM_DPRIIND 59
CPX_PARAM_PRLINEAR 120	CPXPARAM_MIP_Limits_CutsFactor 52	CPX_PARAM_SOLUTIONTARGET (deprecated: see CPXPARAM_OptimalityTarget 106)	CPXPARAM_SolutionType 152	CPX_PARAM_MIPORDIND 90	CPX_PARAM_EACHCUTLIM 60
CPX_PARAM_PREPASS 121	CPXPARAM_MIP_Limits_RampupDetTimeLimit 127	CPXPARAM_SOLUTIONTYPE 152	CPXPARAM_Threads 157	CPX_PARAM_MIPORDTYPE 91	CPX_PARAM_EPAGAP 61
CPX_PARAM_PRESLVND 122	CPXPARAM_MIP_Limits_RampupTimeLimit 128	CPX_PARAM_STARTALG 139	CPXPARAM_TimeLimit 159	CPX_PARAM_MIPSEARCH 92	CPX_PARAM_EPGAP 61
CPX_PARAM_PRICELIM 123	CPXPARAM_MIP_Limits_Solutions 79	CPX_PARAM_STRONGCANDLIM 154	CPXPARAM_Tune_DefTimeLimit 160	CPX_PARAM_MIQCPSTRAT 93	CPX_PARAM_EPINT 62
CPX_PARAM_PROBE 123	CPXPARAM_MIP_Limits_StrongCand 154	CPX_PARAM_STRONGCANDLIM 154	CPXPARAM_Tune_Display 162	CPX_PARAM_MIRCUTS 94	CPX_PARAM_EPMRK 64
CPX_PARAM_PROBEDETTIME 124	CPXPARAM_MIP_Limits_StrongIt 154	CPX_PARAM_STRONGITLIM 154	CPXPARAM_Tune_Measure 163	CPX_PARAM_MPSLONGNUM 94	CPX_PARAM_EPOPT 65
CPX_PARAM_PROBETIME 124	CPXPARAM_MIP_Limits_TreeMemory 160	CPX_PARAM_SUBALG 99	CPXPARAM_Tune_Repeat 164	CPX_PARAM_NETDISPLAY 95	CPX_PARAM_EPPER 65
CPX_PARAM_QPMAKEPSDIND 125	CPXPARAM_MIP_OrderType 91	CPX_PARAM_SUBMIPNODELIMIT 155	CPXPARAM_Tune_TimeLimit 165	CPX_PARAM_NETEPOPT 96	CPX_PARAM_EPRELAX 66
CPX_PARAM_QPMETHOD 138	CPXPARAM_MIP_Pool_AbsGap 146	CPX_PARAM_SYMMETRY 156	CPXPARAM_WorkDir 167	CPX_PARAM_NETEPRHS 96	CPX_PARAM_EPRHS 67
CPX_PARAM_QPNZREADLIM 126	CPXPARAM_MIP_Pool_Capacity 147	CPX_PARAM_THREADS 157	CPXPARAM_WorkMem 168	CPX_PARAM_NETFIND 97	CPX_PARAM_FEASOPTMODE 68
	CPXPARAM_MIP_Pool_Intensity 149	CPX_PARAM_TILIM 159	CraInd 50	CPX_PARAM_NETITLIM 98	CPX_PARAM_FILEENCODING 69
				CPX_PARAM_NETPPRIIND 98	

Algorithm configuration

IP solvers (CPLEX, Gurobi) have a **ton** parameters

- CPLEX has **170-page** manual describing **172** parameters
- Tuning by hand is notoriously **slow, tedious**, and **error-prone**

What's the best **configuration** for the application at hand?



Best configuration for **routing** problems
likely not suited for **scheduling**



How to integrate **machine learning** into **discrete optimization**?

- **Algorithm configuration**
How to tune an algorithm's parameters?
- **Algorithm selection**
Given a variety of algorithms, which to use?
- **Algorithm design**
Can machine learning guide algorithm discovery?

Example: Clustering

Many different algorithms

K-means



Mean shift



Ward



Agglomerative



Birch



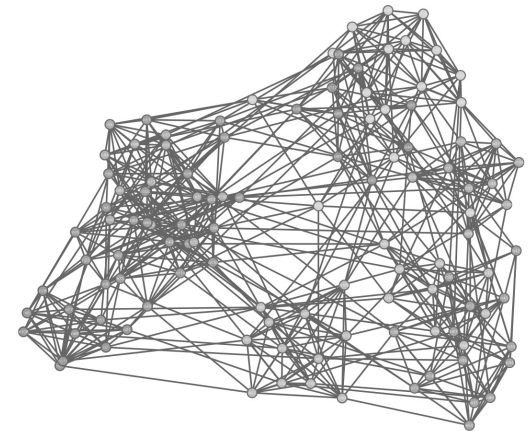
How to **select** the best algorithm for the application at hand?

Algorithm selection in theory

Worst-case analysis has been the main framework for decades
Has led to beautiful, practical algorithms

Worst-case analysis's approach to **algorithm selection**:
Select the algorithm that's best in worst-case scenarios

Worst-case instances **rarely occur in practice**



How to integrate **machine learning** into **discrete optimization**?

- **Algorithm configuration**
How to tune an algorithm's parameters?
- **Algorithm selection**
Given a variety of algorithms, which to use?
- **Algorithm design**
Can machine learning guide algorithm discovery?

How to integrate **machine learning** into **discrete optimization**?

Long-term goal:

Researchers will be empowered with **data-driven tools** to

 Conceive

 Prototype

 Validate

algorithmic ideas...

and provide theoretical guarantees for their discoveries

How to integrate **machine learning** into **discrete optimization**?

Research area is built on a key observation:

**In practice, we have data about
the application domain**

A stack of several cardboard boxes, each secured with a red and white striped string. The boxes are brown and feature a 'FRAGILE' label with three icons: a vertical arrow, a wine glass, and an umbrella. A person's hands are visible at the bottom, holding the stack. The background is blurred, showing a person in a white shirt.

**In practice, we have data about
the application domain**

Routing problems a shipping company solves

**In practice, we have data about
the application domain**



Clustering problems a biology lab solves

**In practice, we have data about
the application domain**



Scheduling problems an airline solves

In practice, we have data about the application domain

How can we use this data to guide:

- **Algorithm configuration**
How to tune an algorithm's parameters?
- **Algorithm selection**
Given a variety of algorithms, which to use?
- **Algorithm design**
Can machine learning guide algorithm discovery?

ML + discrete opt: Potential impact

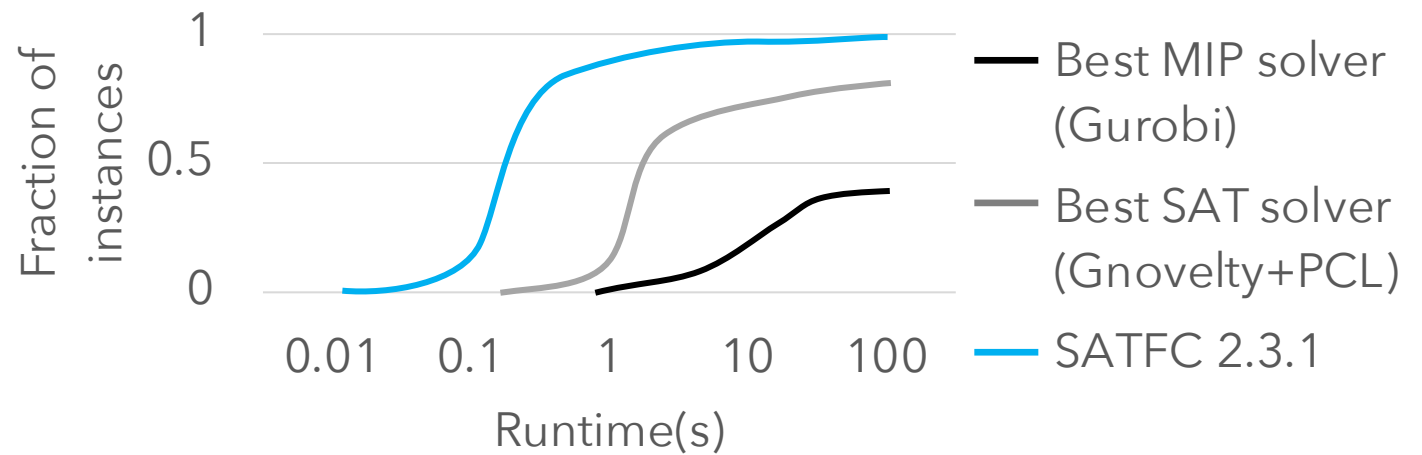
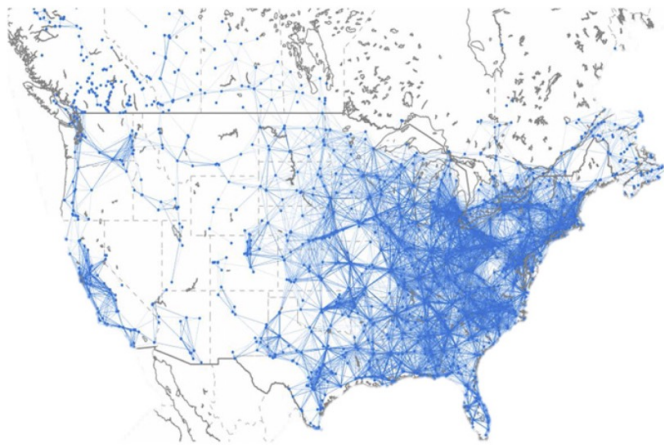
Example: integer programming

- Used heavily throughout industry and science
- **Many** different ways to incorporate **learning** into solving
- Solving is very difficult, so ML can make a huge difference



Example: Spectrum auctions

- In '16-'17, FCC held a \$19.8 billion radio spectrum auction
 - Involves solving huge graph-coloring problems



- SATFC uses algorithm configuration + selection
- Simulations indicate SATFC saved the government billions

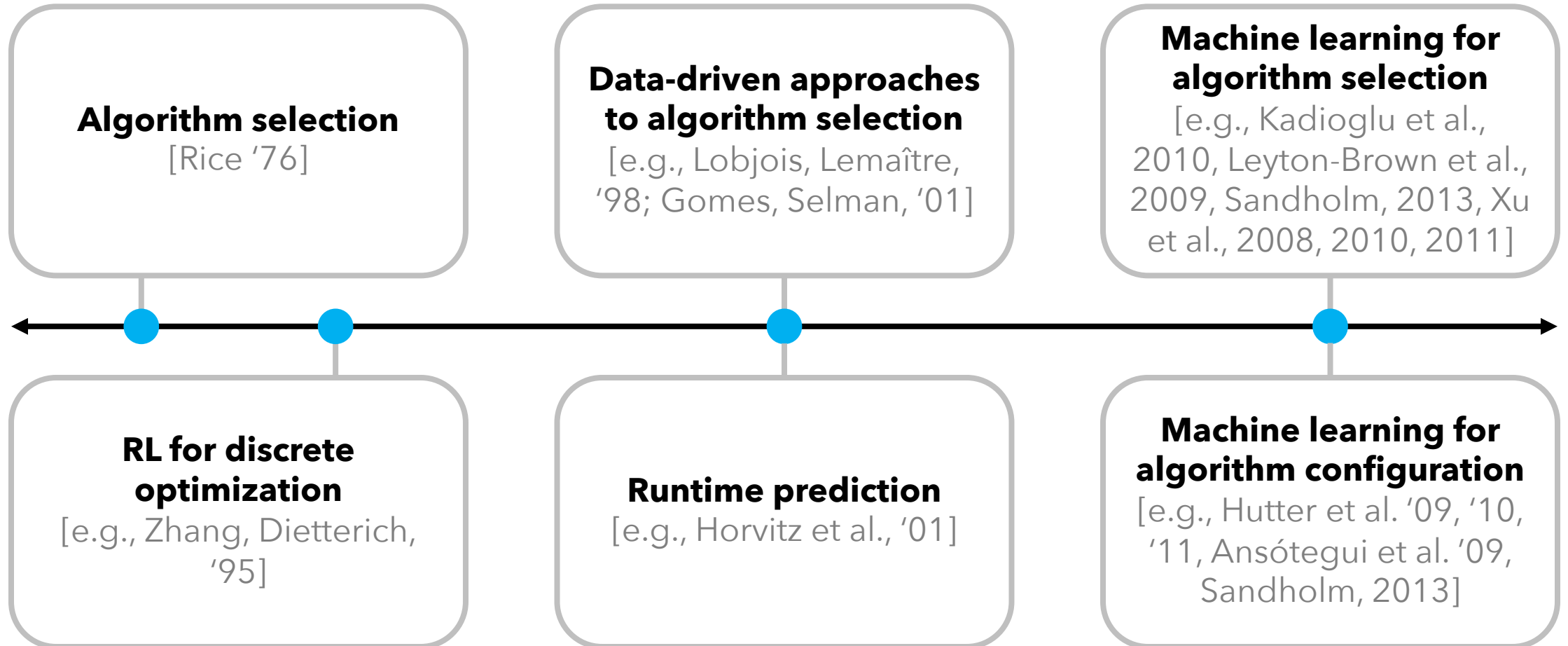
A bit of history

Important research direction in artificial intelligence for decades

Has led to **breakthroughs** in

- Combinatorial auction winner determination
- SAT
- Constraint satisfaction
- Integer programming
- Many other areas

A bit of history



A bit of history



Plan for tutorial

1 Applied techniques

- a. Graph neural networks
- b. Reinforcement learning

2 Theoretical guarantees

- a. Statistical guarantees for algorithm configuration
- b. Algorithms with predictions

Where much of my research has been



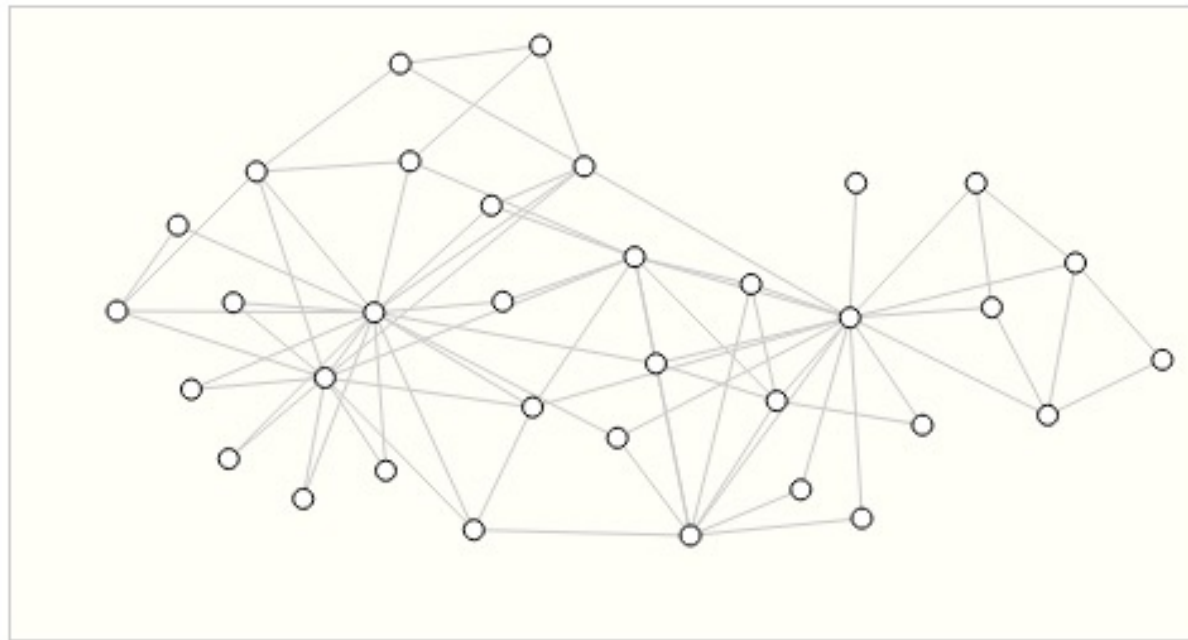
Outline (applied techniques)

- 1. GNNs overview**
2. Integer programming with GNNs
3. Neural algorithmic alignment
4. Learning greedy heuristics with RL

GNN motivation

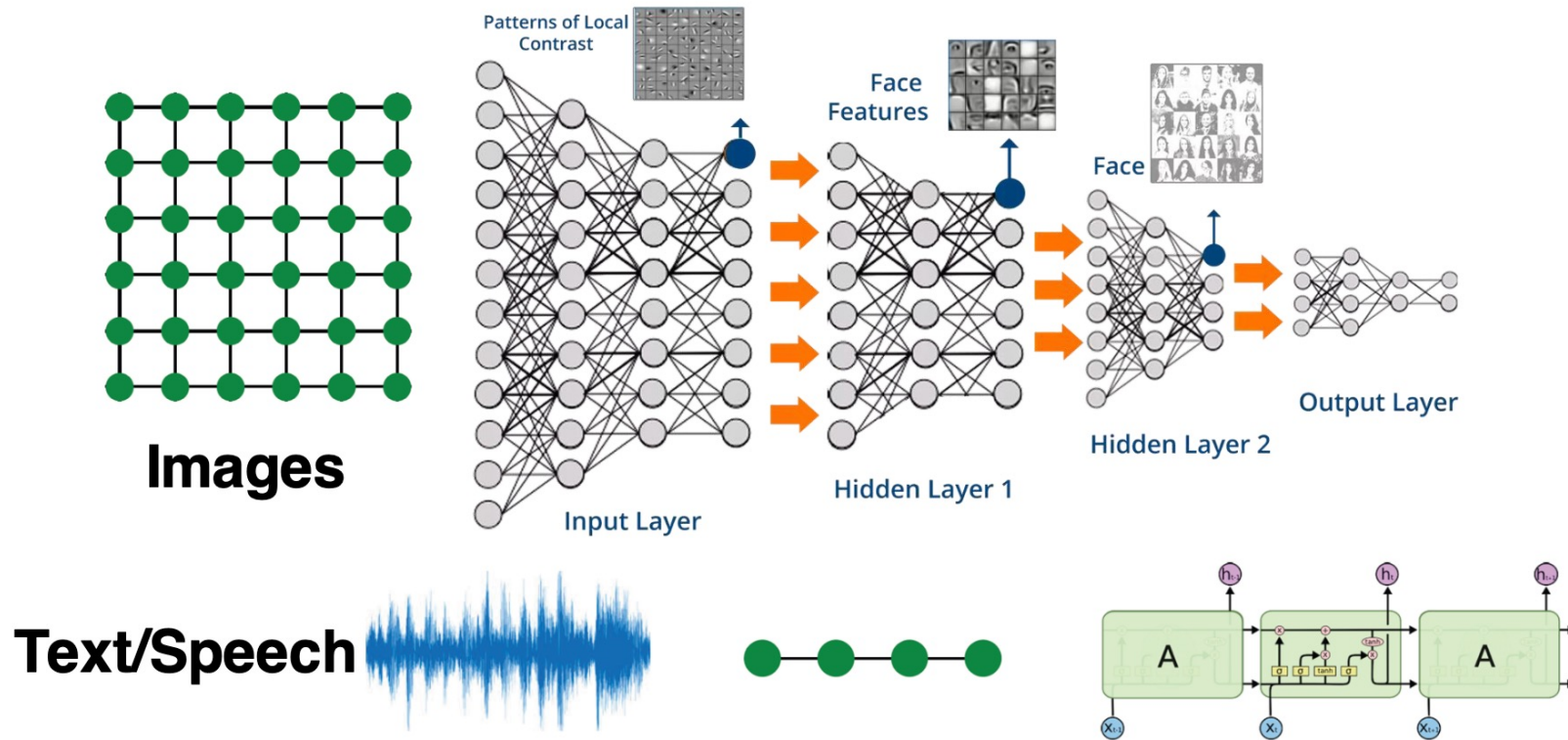
Main question:

How to utilize relational structure for better prediction?



Today: Modern ML toolbox

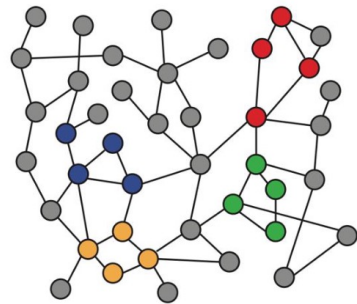
Modern DL toolbox is designed for simple sequences & grids



Why is graph deep learning hard?

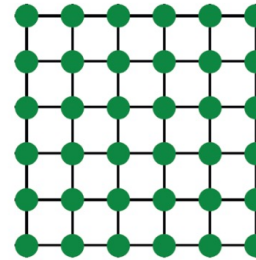
Networks are complex

- Arbitrary size and complex topological structure



Networks

versus



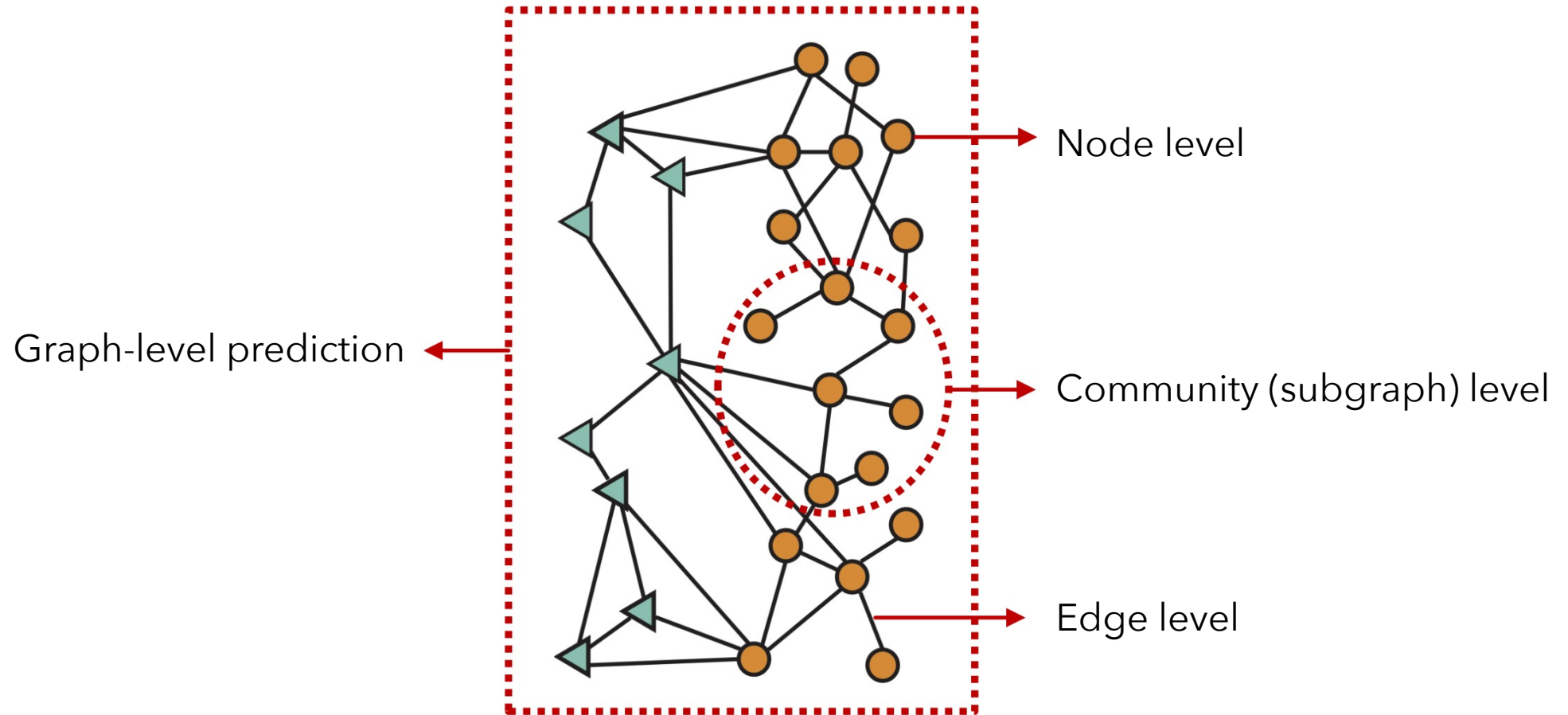
Images



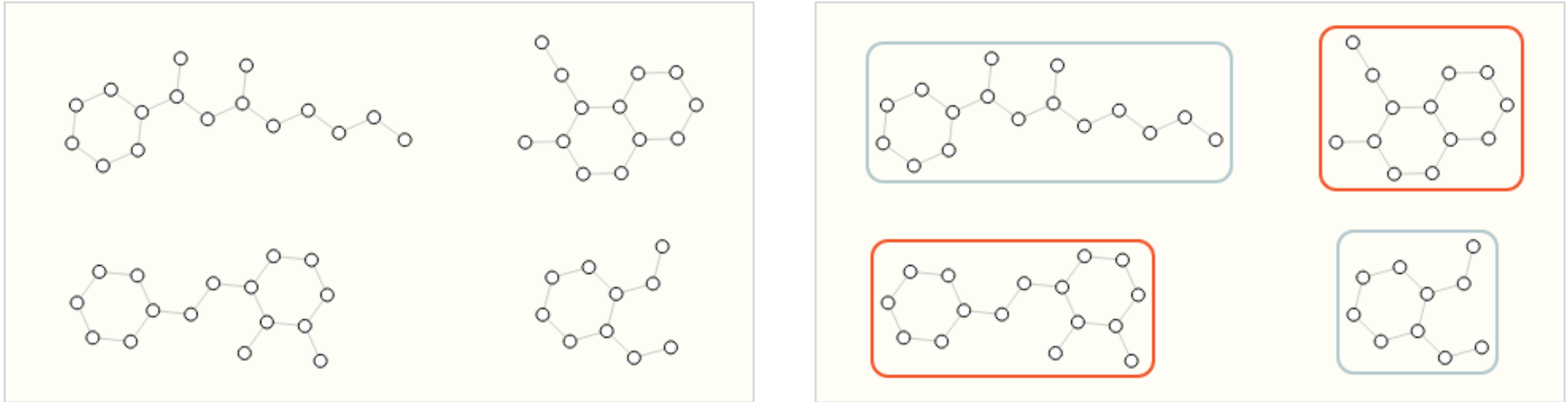
Text

- No fixed node ordering or reference point
- Often dynamic and have multimodal features

Different types of tasks



Prediction with graphs: Examples

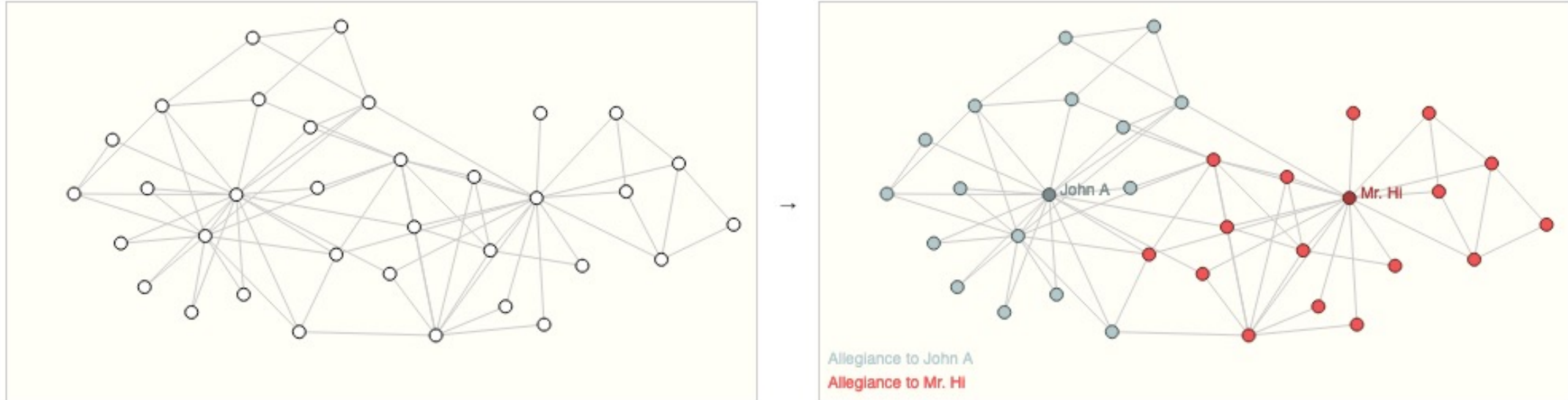


Graph-level tasks:

E.g., for a molecule represented as a graph, could predict:

- What the molecule smells like
- Whether it will bind to a receptor implicated in a disease

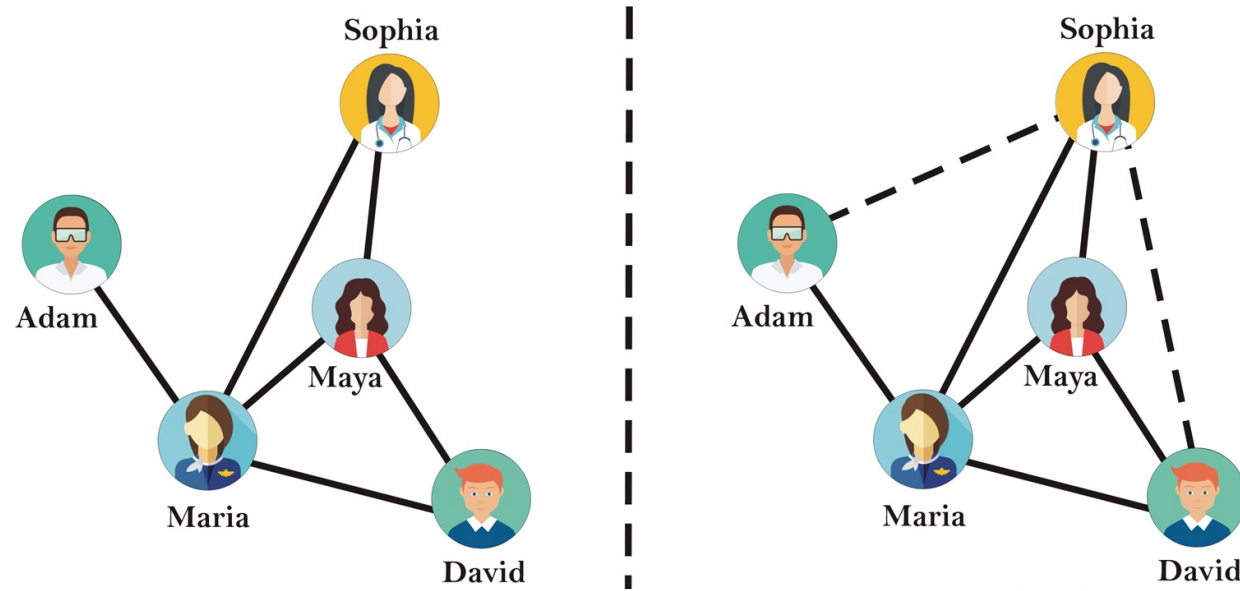
Prediction with graphs: Examples



Node-level tasks:

E.g., political affiliations of users in a social network

Prediction with graphs: Examples



Edge-level tasks: E.g.:

- Suggesting new friends
- Recommendations on Amazon, Netflix, ...

Example: Traffic routing



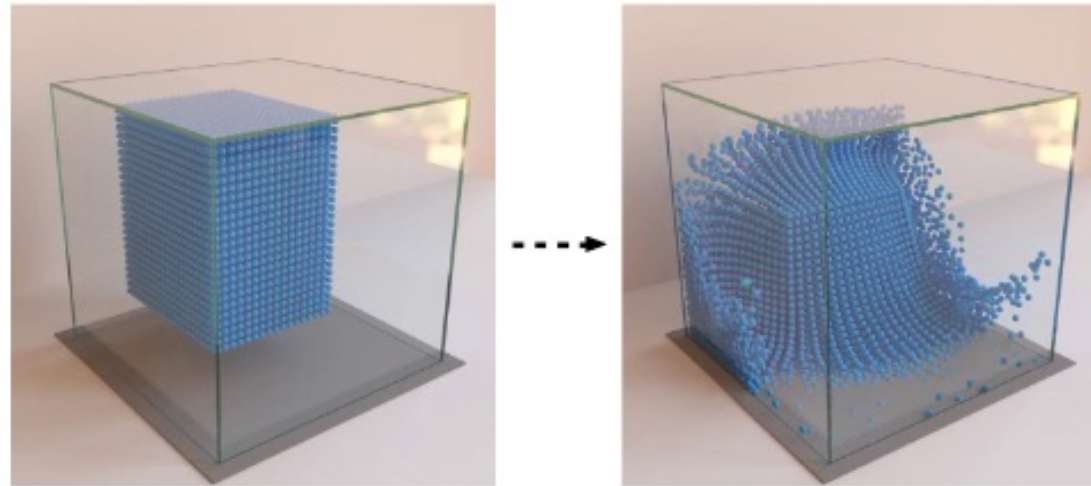
E.g., Google maps

deepmind.com/blog/article/traffic-prediction-with-advanced-graph-neural-networks

Example: Learning to simulate physics

Nodes: Particles

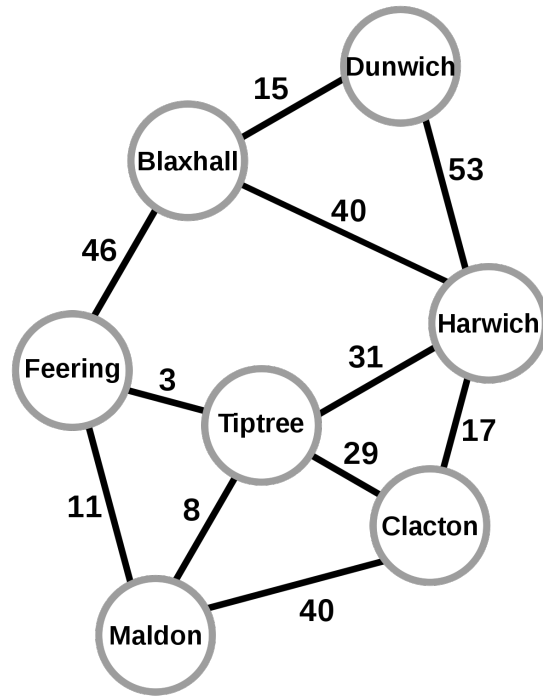
Edges: Interaction between particles



Goal: Predict how a graph will evolve over time

Example: Combinatorial optimization

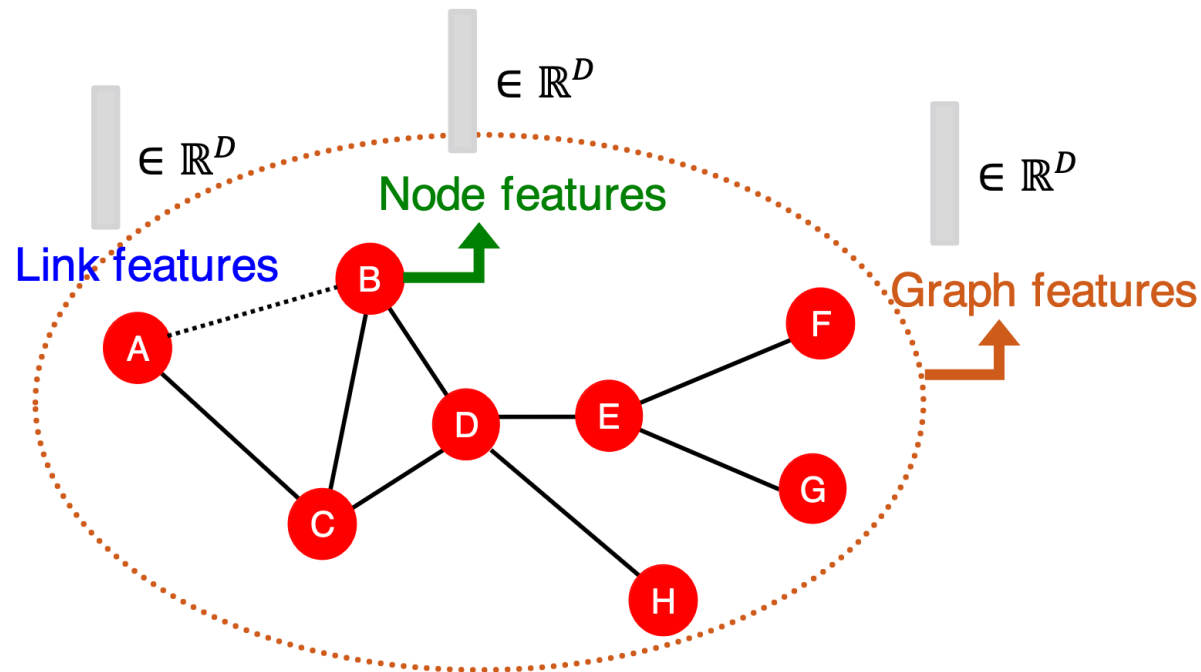
Replace full algorithm or learn steps (e.g., branching decision)



$$\begin{array}{ll} \text{maximize} & \mathbf{c} \cdot \mathbf{z} \\ \text{subject to} & \mathbf{A}\mathbf{z} \leq \mathbf{b} \\ & \mathbf{z} \in \mathbb{Z}^n \end{array}$$

Graph neural networks: First step

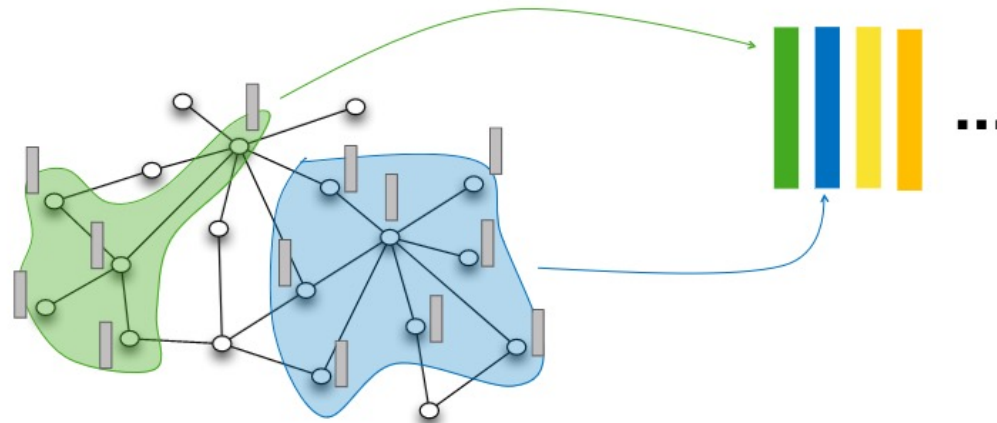
- Design features for nodes/links/graphs
- Obtain features for all training data



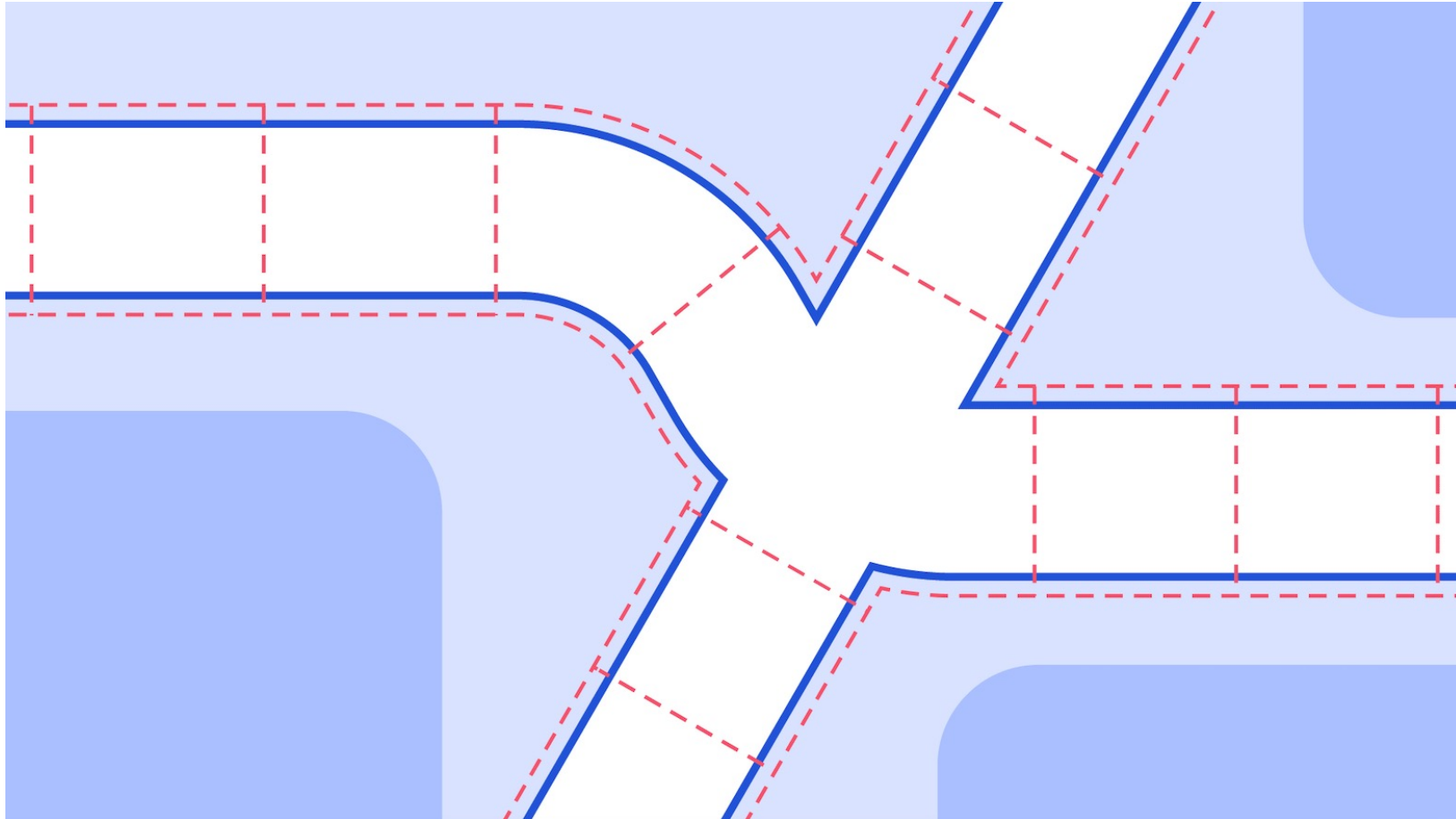
Graph neural networks: Objective

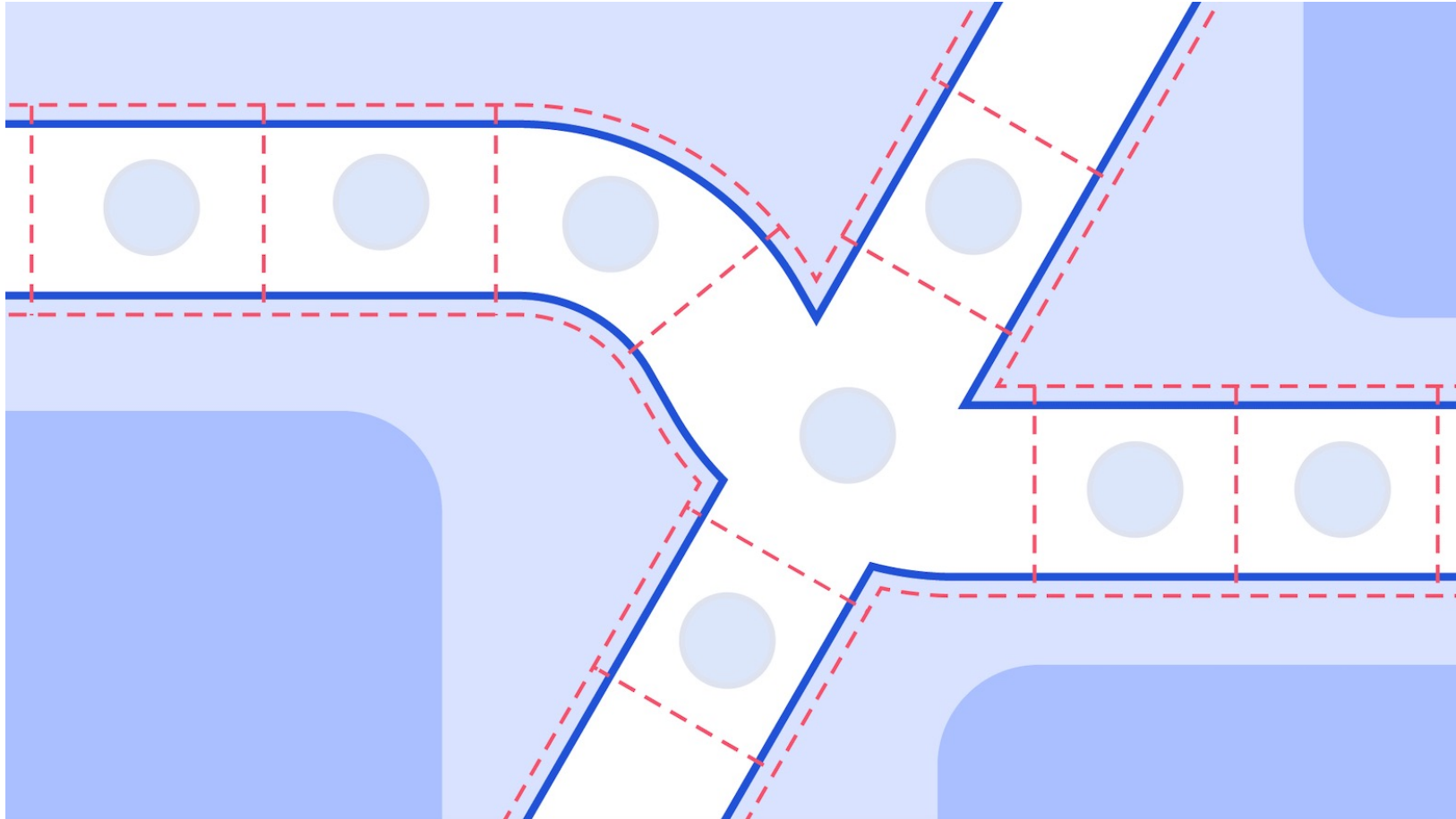
Idea:

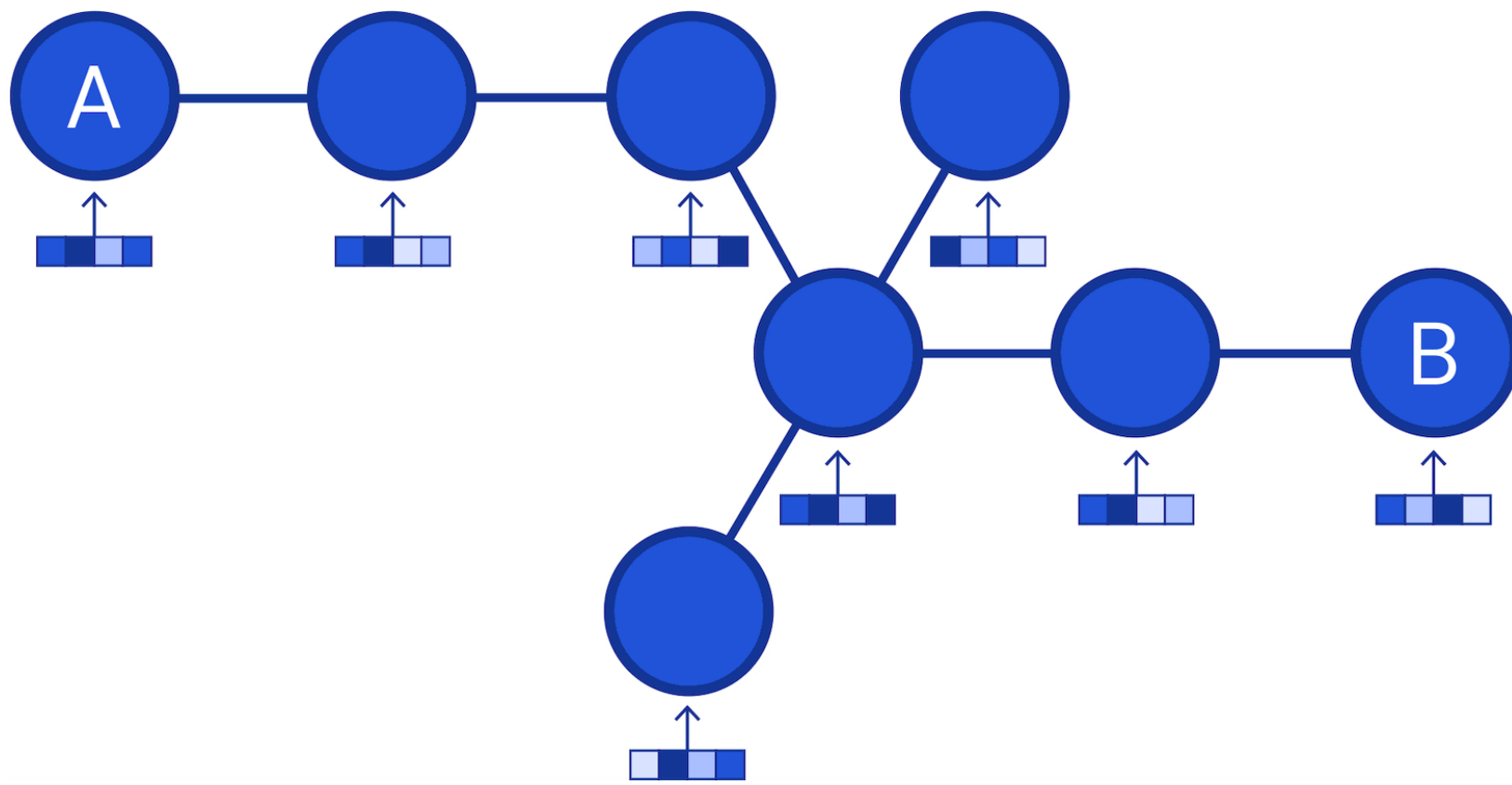
1. Encode each node and its neighborhood with embedding
2. Aggregate set of node embeddings into graph embedding
3. Use embeddings to make predictions

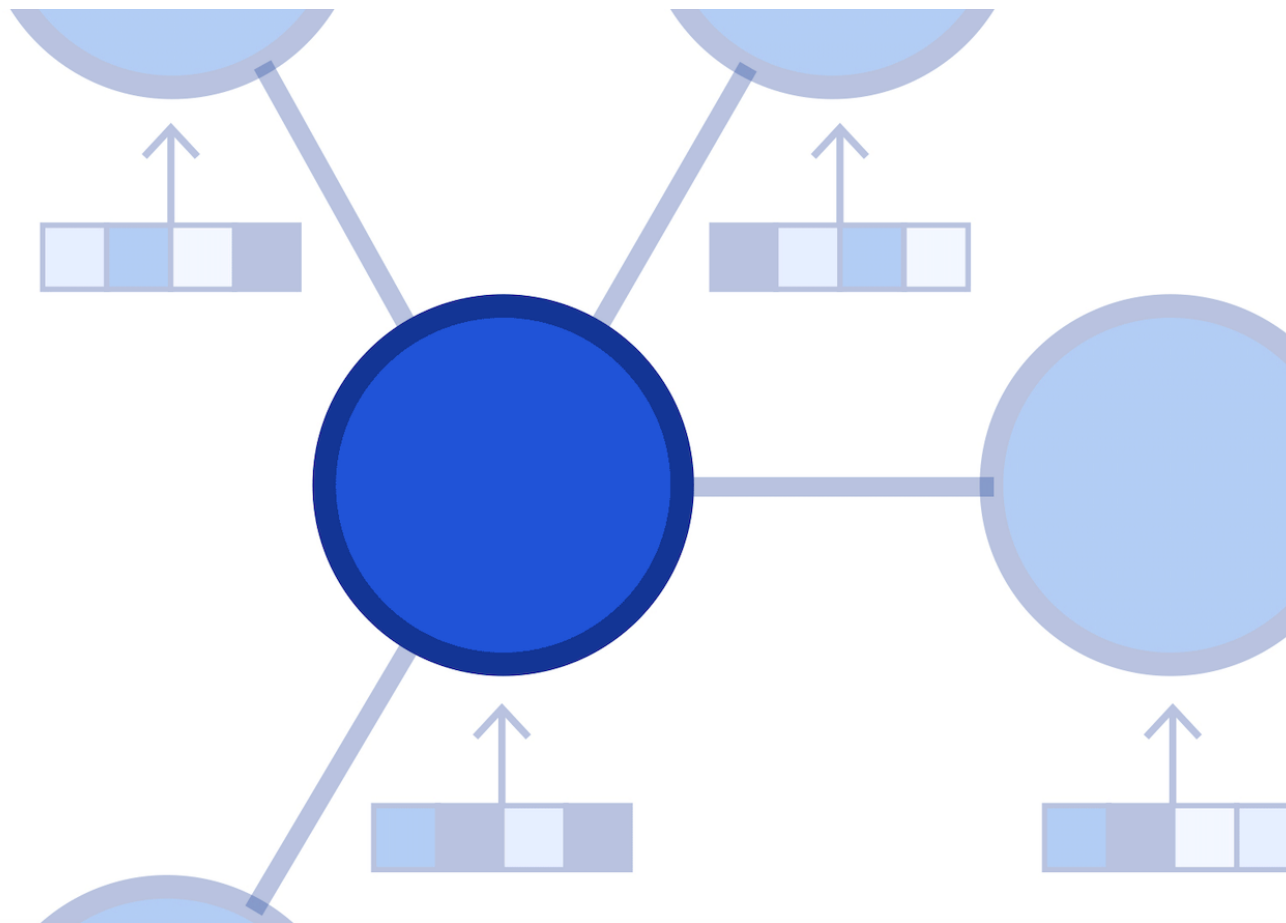


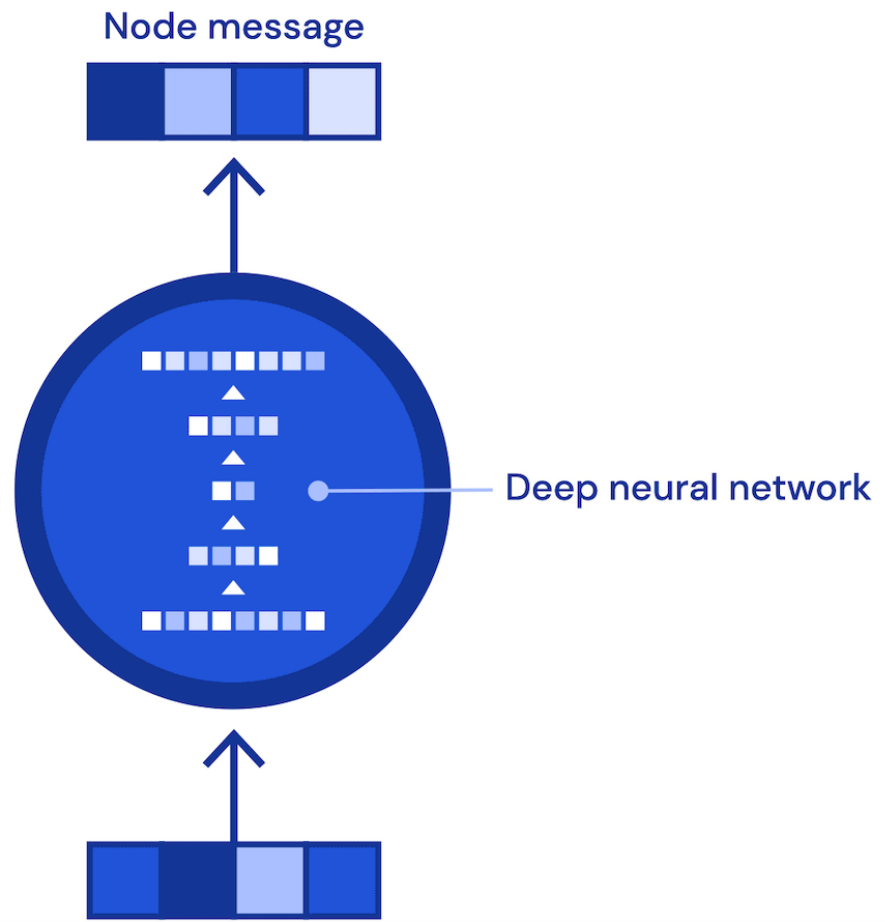


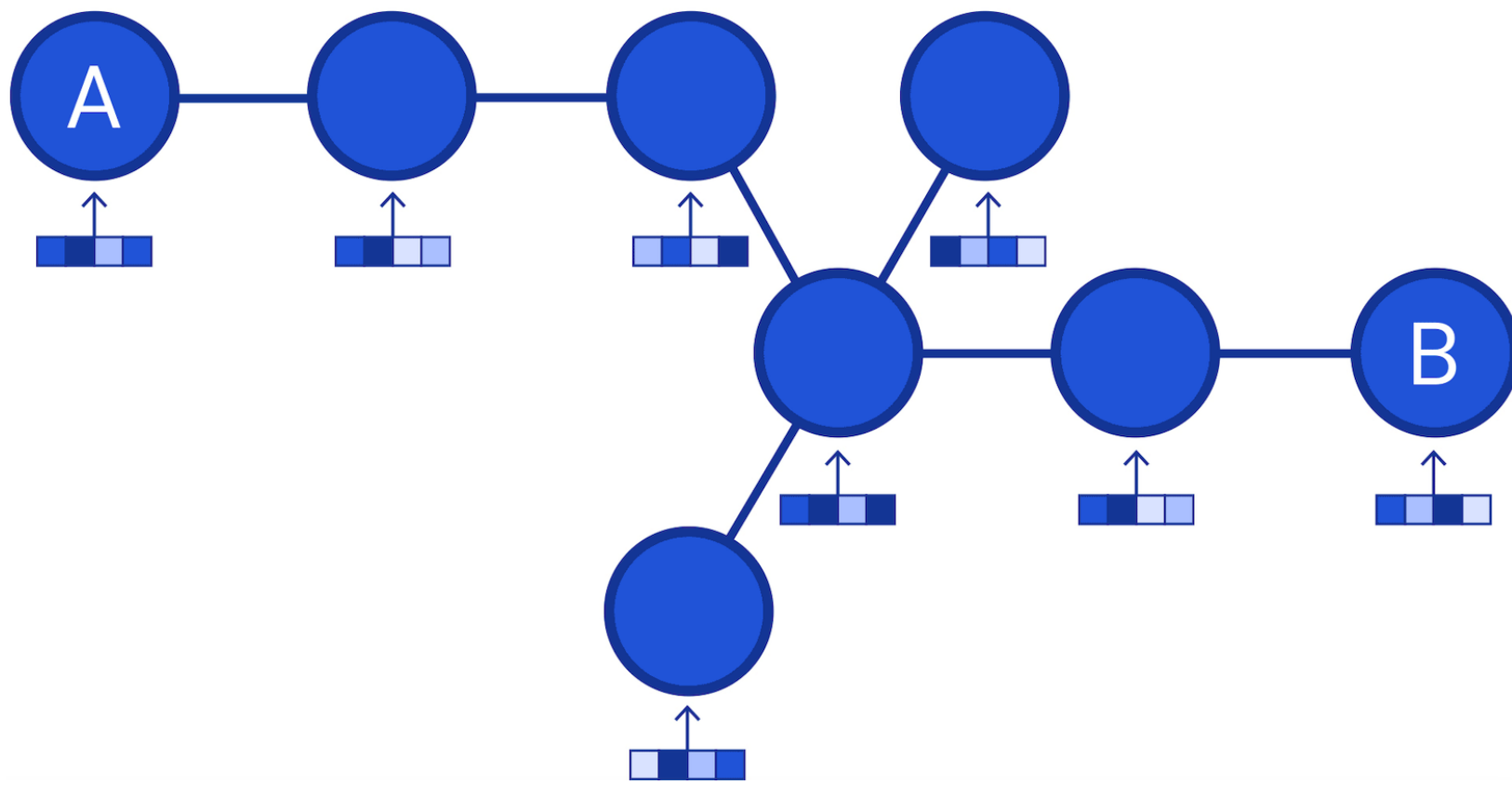


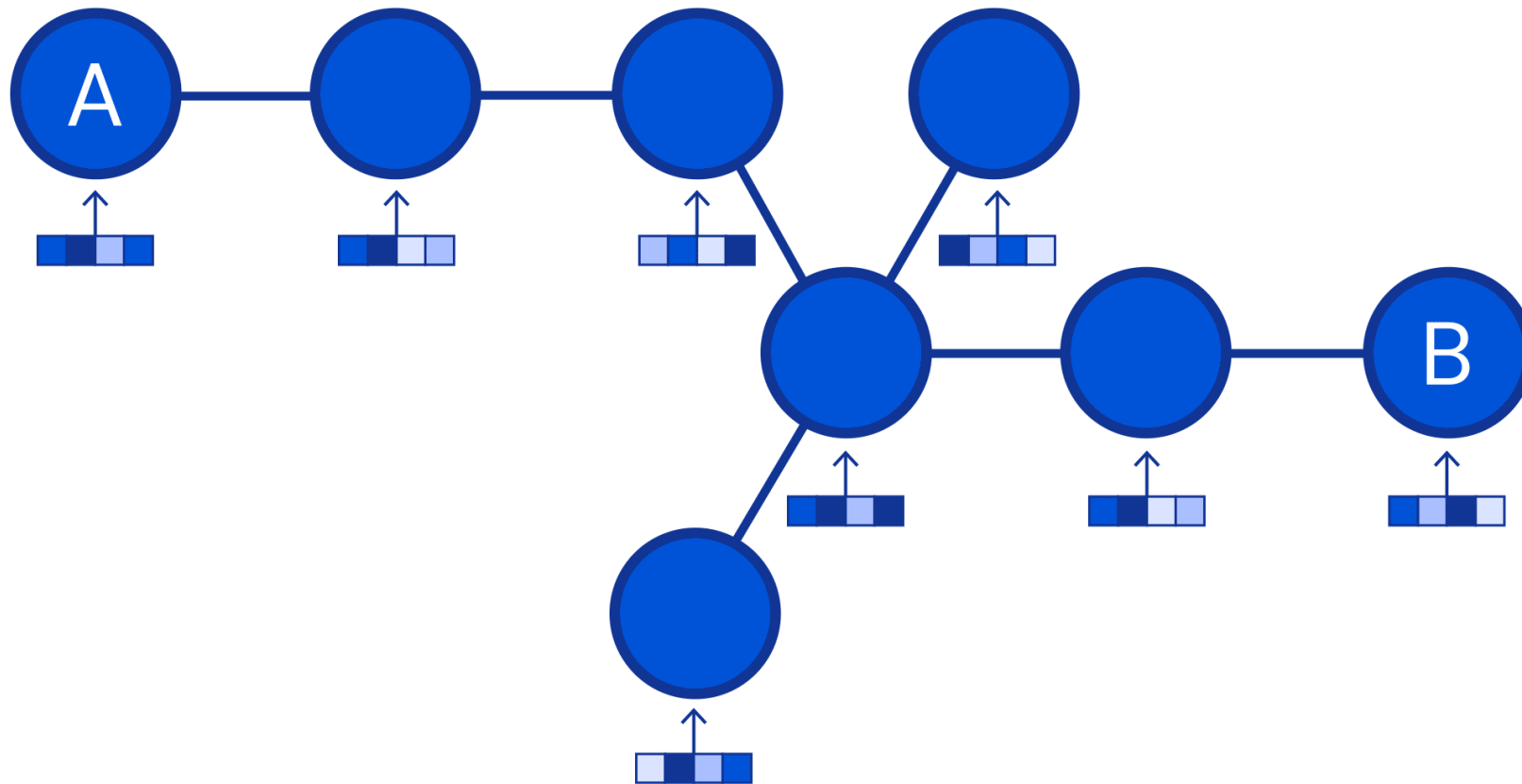












Encoding neighborhoods: General form

$\mathbf{h}_v^{(0)} = \mathbf{x}_v$ (feature representation for node v)

In each round $k \in [K]$, for each node v :

1. **Aggregate** over neighbors

$$\mathbf{m}_{N(v)}^{(k)} = \text{AGGREGATE}^{(k)} \left(\left\{ \mathbf{h}_u^{(k-1)} : u \in N(v) \right\} \right)$$

Neighborhood of v

Encoding neighborhoods: General form

$\mathbf{h}_v^{(0)} = \mathbf{x}_v$ (feature representation for node v)

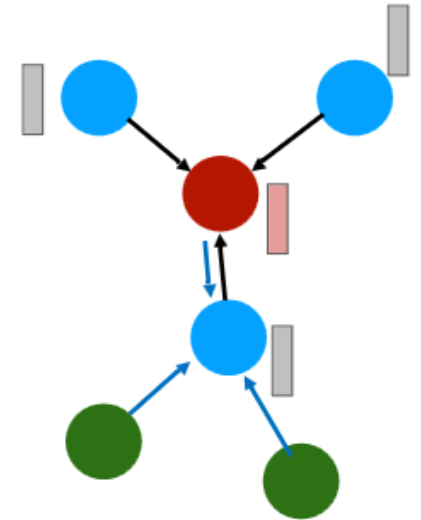
In each round $k \in [K]$, for each node v :

1. **Aggregate** over neighbors

$$\mathbf{m}_{N(v)}^{(k)} = \text{AGGREGATE}^{(k)} \left(\left\{ \mathbf{h}_u^{(k-1)} : u \in N(v) \right\} \right)$$

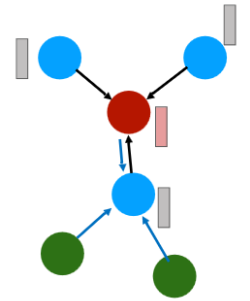
2. **Update** current node representation

$$\mathbf{h}_v^{(k)} = \text{COMBINE}^{(k)} \left(\mathbf{h}_v^{(k-1)}, \mathbf{m}_{N(v)}^{(k)} \right)$$



The basic GNN

[Merkwirth and Lengauer '05; Scarselli et al. '09]



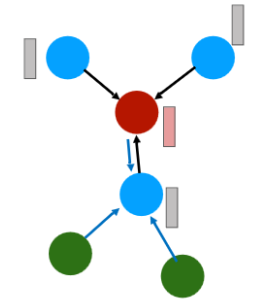
$$\mathbf{m}_{N(v)} = \text{AGGREGATE}(\{\mathbf{h}_u : u \in N(v)\}) = \sum_{u \in N(v)} \mathbf{h}_u$$

$$\text{COMBINE}(\mathbf{h}_v, \mathbf{m}_{N(v)}) = \sigma(\underbrace{W_{\text{self}}}_{\text{Trainable parameters}} \mathbf{h}_v + \underbrace{W_{\text{neigh}}}_{\text{Trainable parameters}} \mathbf{m}_{N(v)} + \underbrace{\mathbf{b}}_{\text{Trainable parameters}})$$

Non-linearity (e.g.,
tanh or ReLU)

Trainable parameters

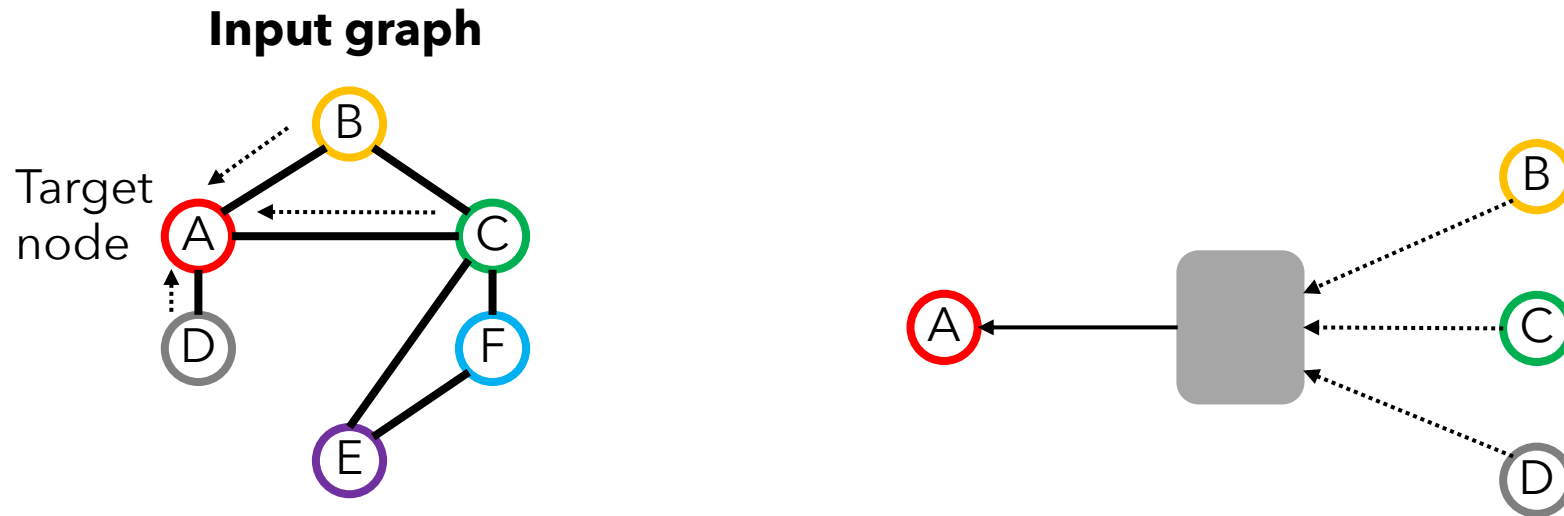
Aggregation functions



$$\mathbf{m}_{N(v)} = \text{AGGREGATE}(\{\mathbf{h}_u : u \in N(v)\}) = \bigoplus_{u \in N(v)} \mathbf{h}_u$$

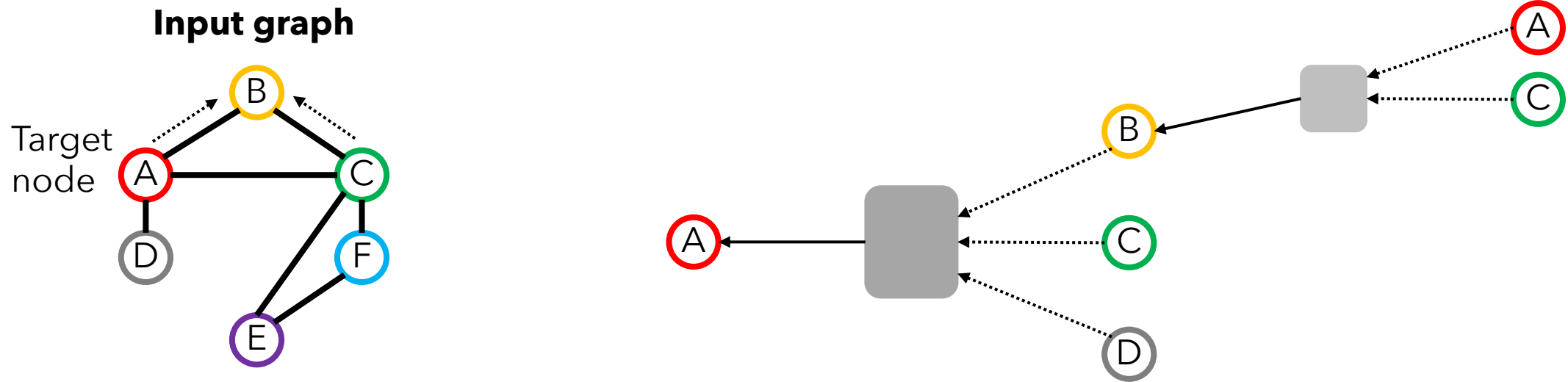
Other element-wise aggregators, e.g.:
Maximization, averaging

Node embeddings unrolled



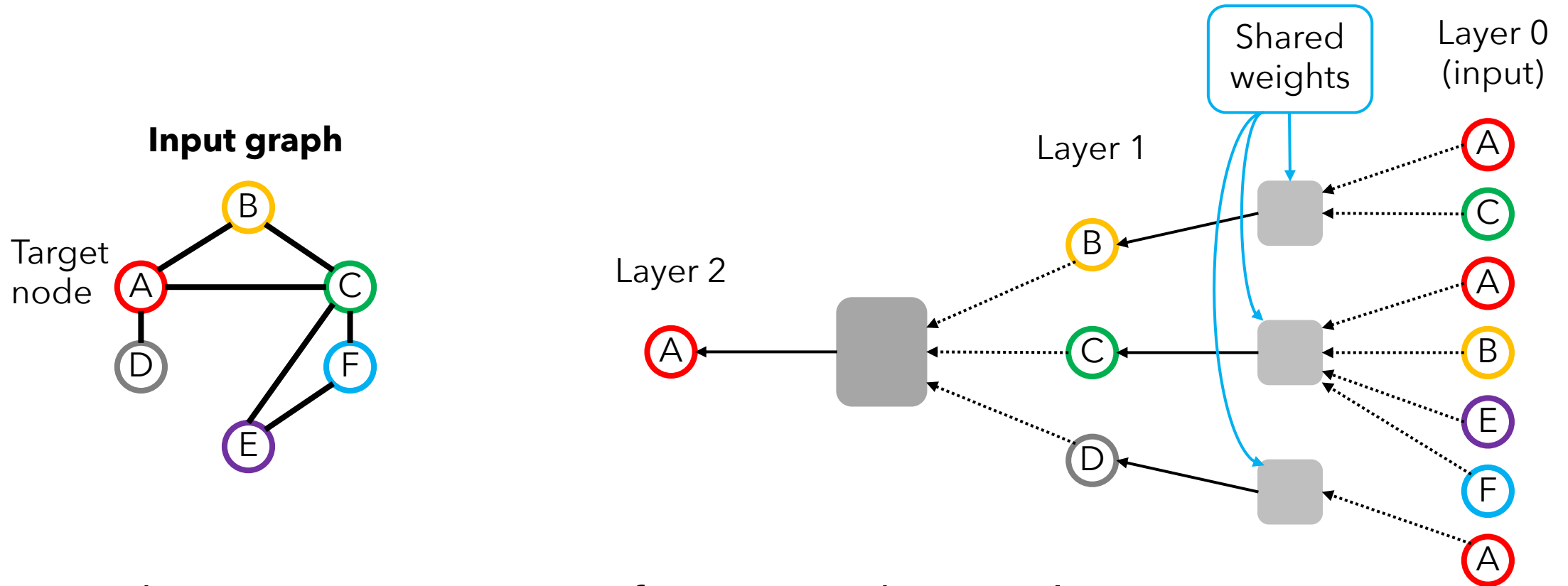
Grey boxes: aggregation functions that we learn

Node embeddings unrolled



Grey boxes: aggregation functions that we learn

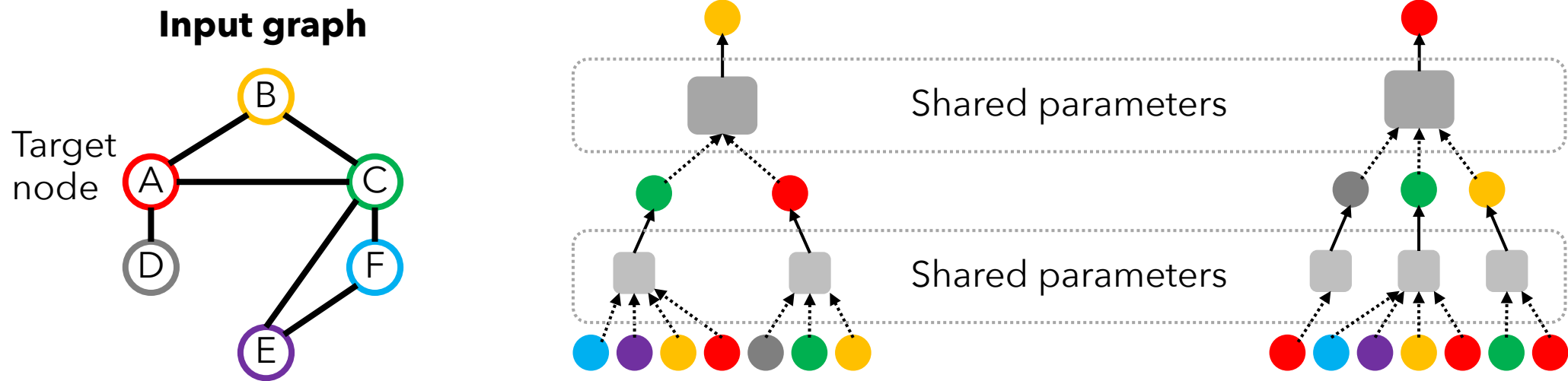
Node embeddings unrolled



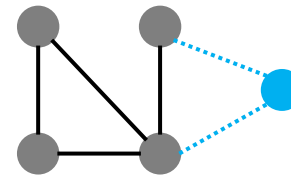
Grey boxes: aggregation functions that we learn

Node embeddings unrolled

Use the same aggregation functions for all nodes



Can generate encodings for previously unseen nodes & graphs!



Outline (applied techniques)

1. GNNs overview
- 2. Integer programming with GNNs**
3. Neural algorithmic alignment
4. Learning greedy heuristics with RL

Gasse, Chételat, Ferroni, Charlin, Lodi; NeurIPS'19

Integer programming solvers

Most popular tool for solving combinatorial problems



Routing



Manufacturing



Scheduling



Planning



Finance

Integer and linear programming

Integer program (IP)

$$\begin{aligned} \max \quad & \mathbf{c} \cdot \mathbf{z} \\ \text{s.t.} \quad & A\mathbf{z} \leq \mathbf{b} \\ & \mathbf{z} \in \mathbb{Z}^n \end{aligned}$$

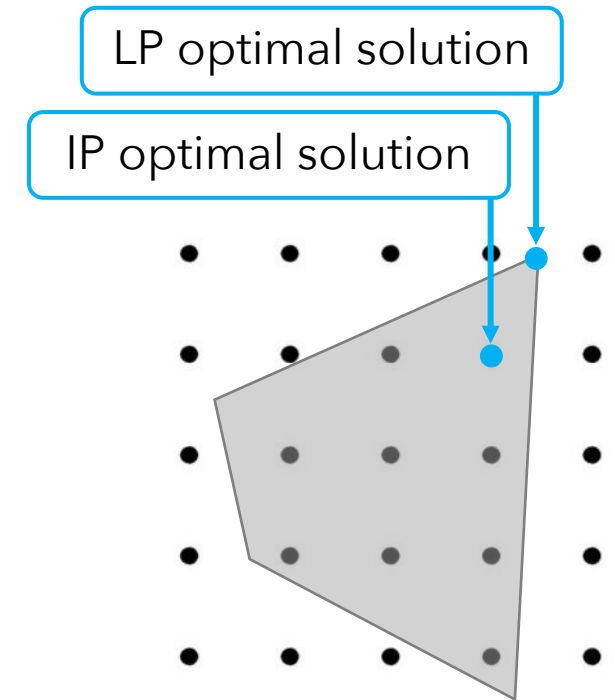
NP-hard

Linear program (LP)

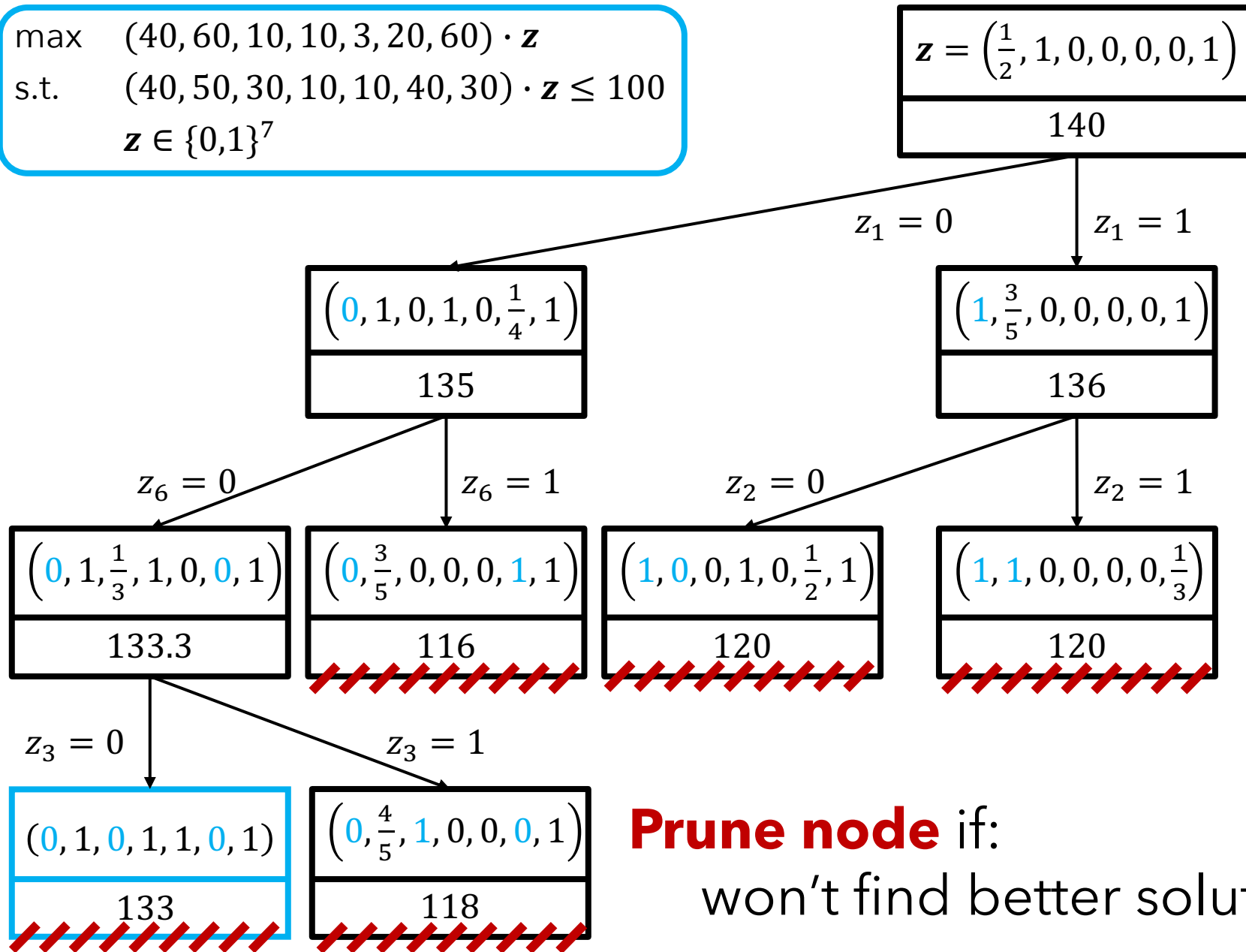
$$\begin{aligned} \max \quad & \mathbf{c} \cdot \mathbf{z} \\ \text{s.t.} \quad & A\mathbf{z} \leq \mathbf{b} \\ & \mathbf{z} \in \mathbb{Z}^n \end{aligned}$$

Efficiently solvable

LP provides valuable guidance in B&B



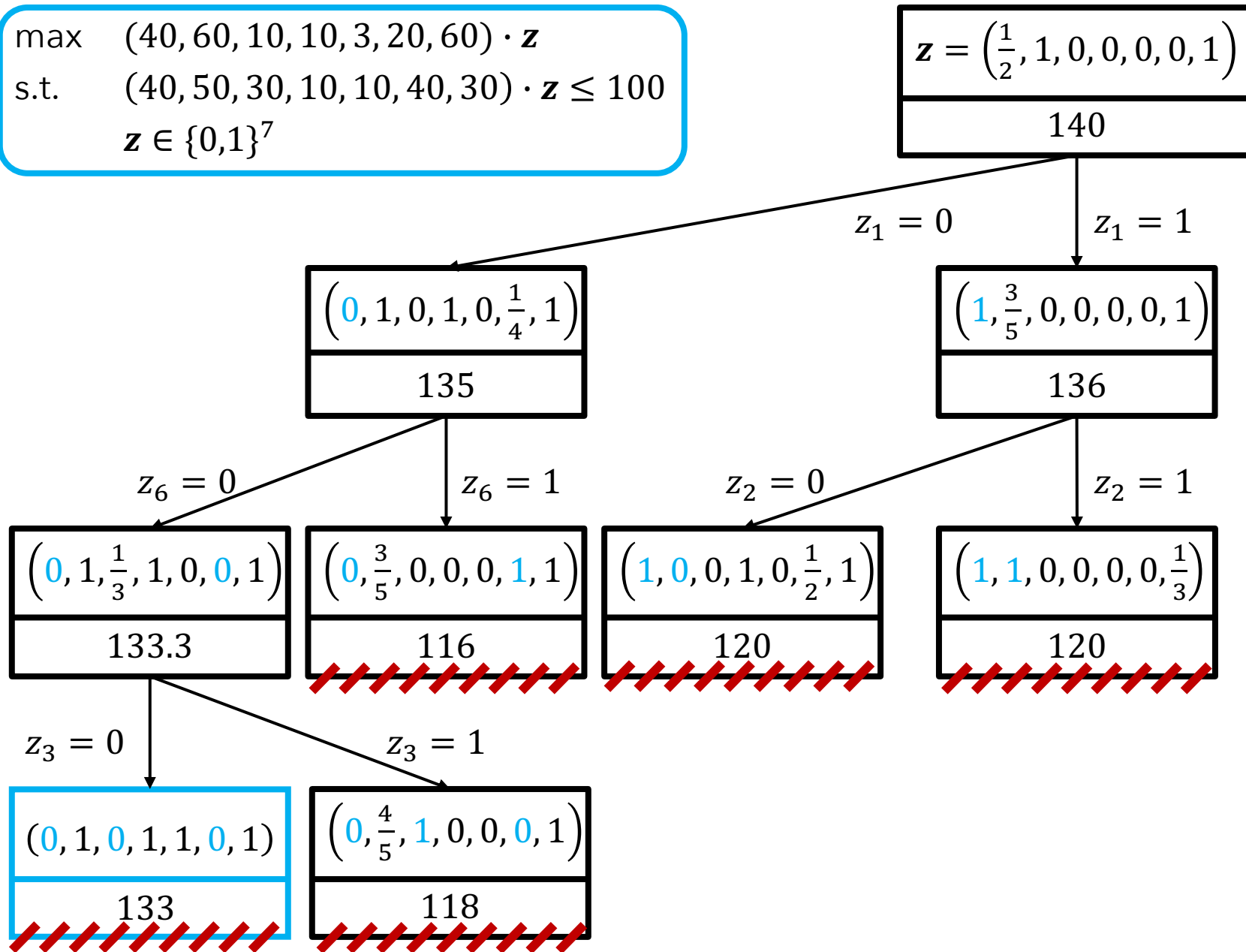
$$\begin{aligned} \max \quad & (40, 60, 10, 10, 3, 20, 60) \cdot \mathbf{z} \\ \text{s.t.} \quad & (40, 50, 30, 10, 10, 40, 30) \cdot \mathbf{z} \leq 100 \\ & \mathbf{z} \in \{0,1\}^7 \end{aligned}$$



Branch
and
bound
(B&B)

Prune node if:
won't find better solution along branch

$$\begin{aligned} \max \quad & (40, 60, 10, 10, 3, 20, 60) \cdot \mathbf{z} \\ \text{s.t.} \quad & (40, 50, 30, 10, 10, 40, 30) \cdot \mathbf{z} \leq 100 \\ & \mathbf{z} \in \{0,1\}^7 \end{aligned}$$



This section:
Variable selection

Variable selection policies (VSPs)

Score-based variable selection policies:

At leaf Q , branch on variable z_i maximizing **score**(Q, i) $\in \mathbb{R}$

Many options! Little known about which to use when

Gauthier, Ribière, Math. Prog. '77; Beale, Annals of Discrete Math. '79; Linderoth, Savelsbergh, INFORMS JoC '99; Achterberg, Math. Prog. Computation '09; Gilpin, Sandholm, Disc. Opt. '11; ...

Variable selection policy example

At node j with LP objective value $z(j)$:

- Let $z_i^+(j)$ be the LP objective value after setting $x_i = 1$
- Let $z_i^-(j)$ be the LP objective value after setting $x_i = 0$

VSP example: Branch on the variable x_i that maximizes
$$(z(j) - z_i^+(j))(z(j) - z_i^-(j))$$

In more detail, scoring rule is $\max\{z(j) - z_i^+(j), 10^{-6}\} \cdot \max\{z(j) - z_i^-(j), 10^{-6}\}$:

If $z(j) - z_i^+(j) = 0$, would lose information stored in $z(j) - z_i^-(j)$

Strong branching

Challenge: Computing $z_i^-(j), z_i^+(j)$ requires solving a lot of LPs

- Computing all LP relaxations referred to as **strong-branching**
- Very **time intensive**

Pro: Strong branching leads to small search trees

Idea: Train an ML model to imitate strong-branching

Khalil et al. [AAAI'16], Alvarez et al. [INFORMS JoC'17], Hansknecht et al. [arXiv'18]

This section: using a GNN

Outline (applied techniques)

1. GNNs overview
2. Integer programming with GNNs
 - i. Machine learning formulation**
 - ii. Baselines
 - iii. Experiments
 - iv. Additional research
3. Neural algorithmic alignment
4. Learning greedy heuristics with RL

Problem formulation

Goal: learn a policy $\pi(x_i | s_t)$

Probability of branching on variable x_i when solver is in state s_t

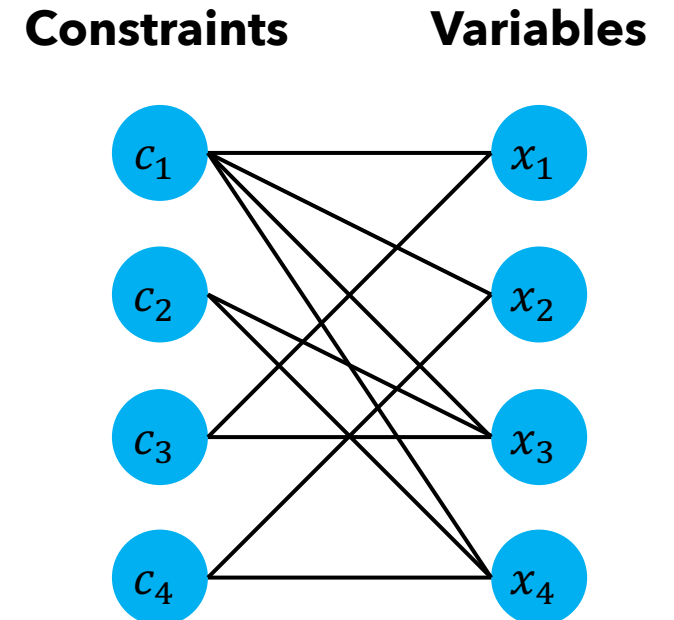
Approach (imitation learning):

- Run strong branching on training set of instances
- Collect dataset of (state, variable) pairs $S = \{(s_i, x_{i^*})\}_{i=1}^N$
- Learn policy π_θ with training set S

State encoding

State s_t of B&B encoded as a **bipartite graph**
with **node** and **edge features**

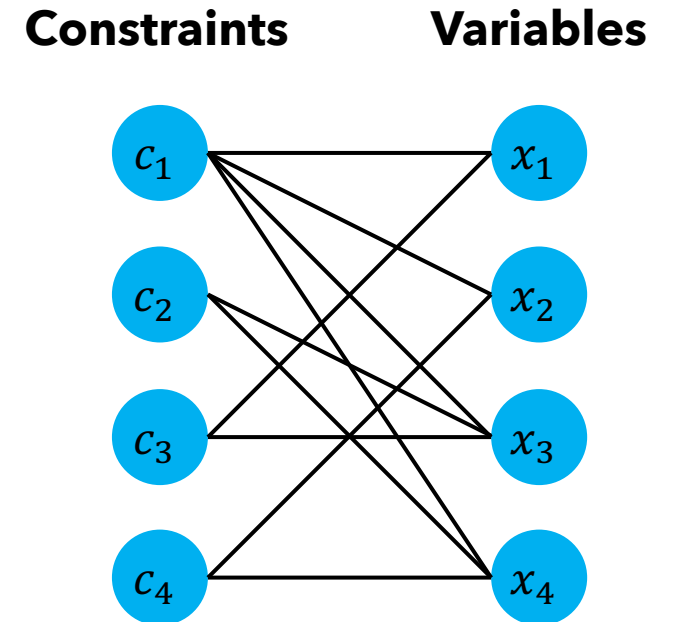
$$\begin{aligned} \max \quad & 9x_1 + 5x_2 + 6x_3 + 4x_4 \\ \text{s.t.} \quad & 6x_1 + 3x_2 + 5x_3 + 2x_4 \leq 10 \quad (c_1) \\ & x_3 + x_4 \leq 10 \quad (c_2) \\ & -x_1 + x_3 \leq 0 \quad (c_3) \\ & -x_2 + x_4 \leq 0 \quad (c_4) \\ & x_1, x_2, x_3, x_4 \in \{0,1\} \end{aligned}$$



State encoding

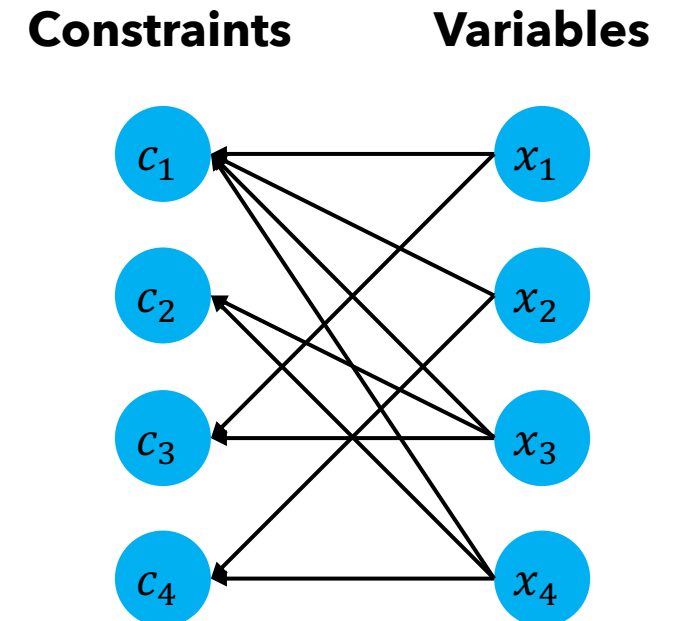
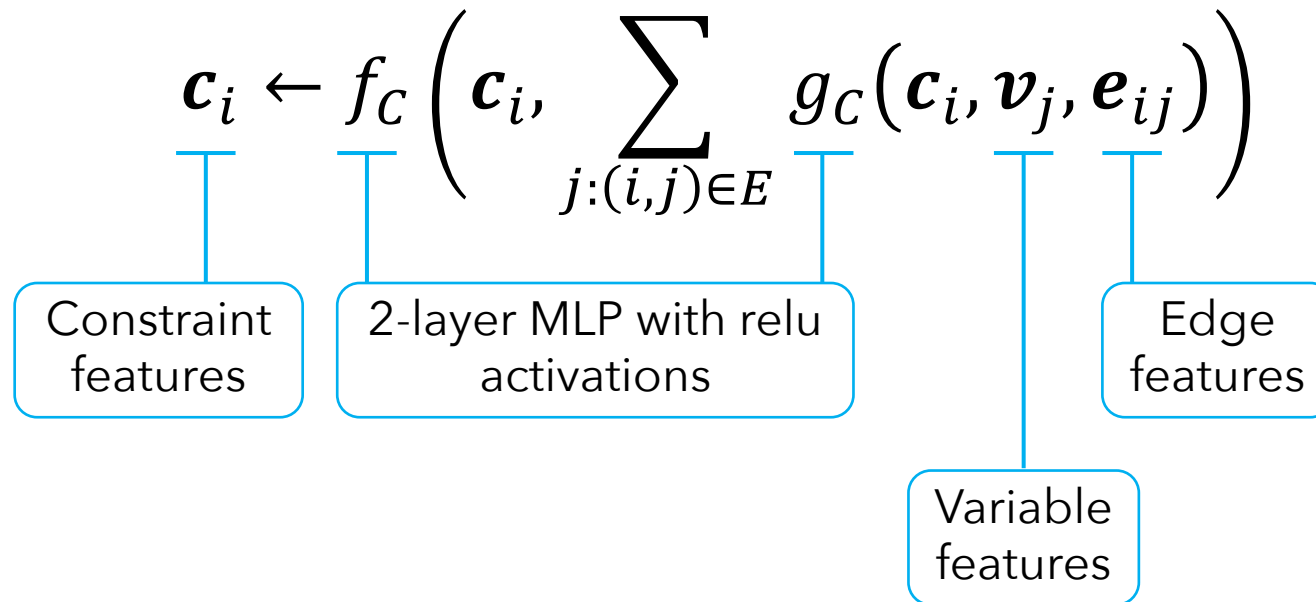
State s_t of B&B encoded as a **bipartite graph** with **node** and **edge features**

- **Edge feature:** constraint coefficient
- **Example node features:**
 - Constraints:
 - Cosine similarity with objective
 - Tight in LP solution?
 - Variables:
 - Objective coefficient
 - Solution value equals upper/lower bound?



GNN structure

1. Pass from variables \rightarrow constraints



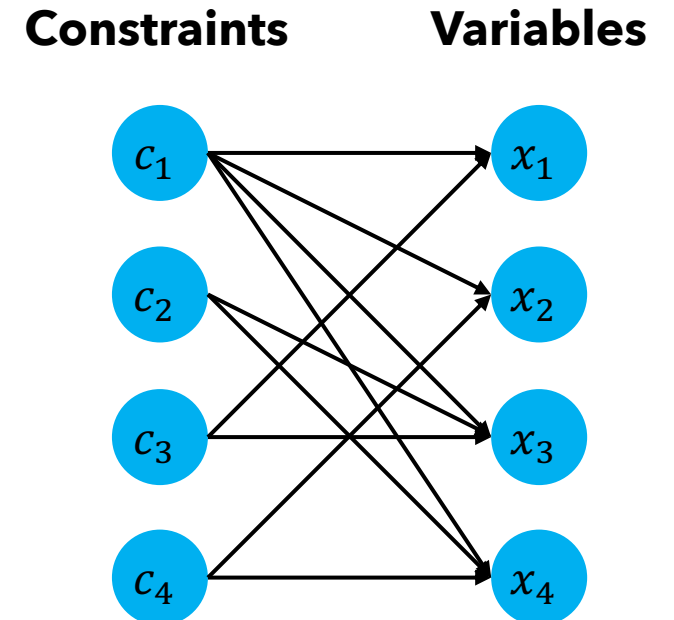
GNN structure

1. Pass from variables \rightarrow constraints

$$\mathbf{c}_i \leftarrow f_C \left(\mathbf{c}_i, \sum_{j:(i,j) \in E} g_C(\mathbf{c}_i, \mathbf{v}_j, \mathbf{e}_{ij}) \right)$$

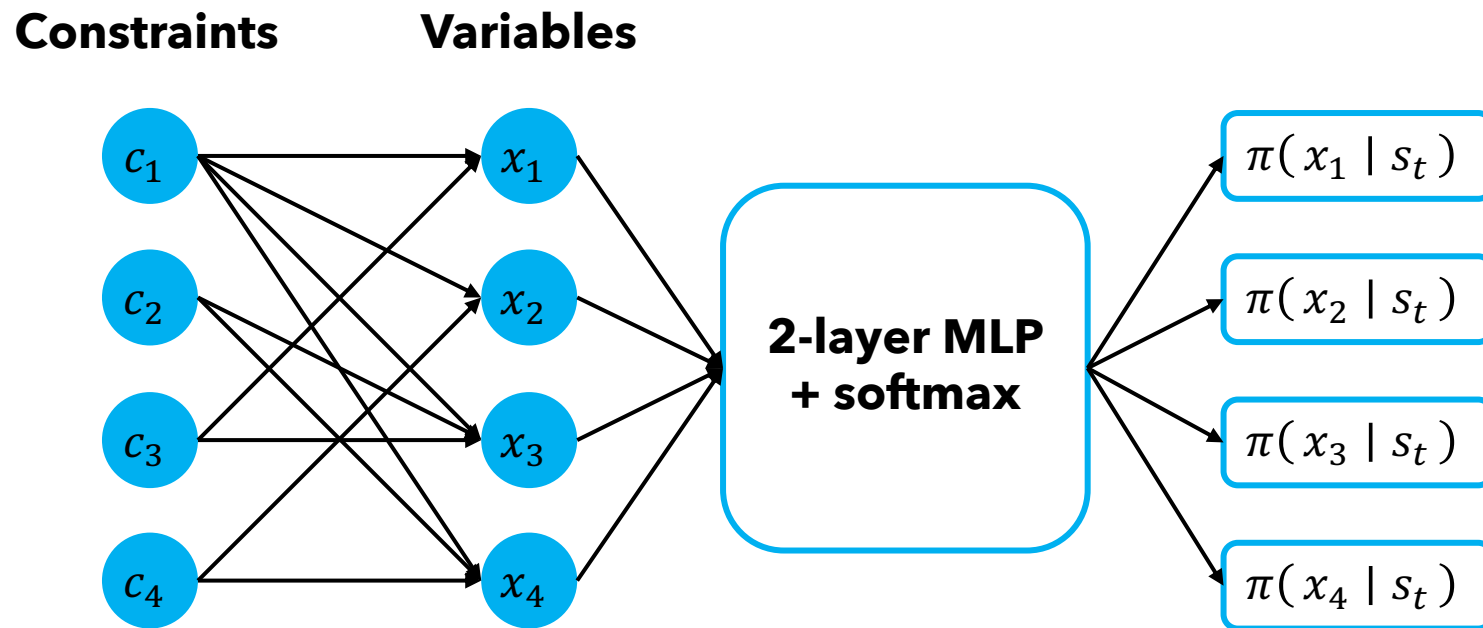
2. Pass from constraints \rightarrow variables

$$\mathbf{v}_j \leftarrow f_V \left(\mathbf{v}_j, \sum_{i:(i,j) \in E} g_V(\mathbf{c}_i, \mathbf{v}_j, \mathbf{e}_{ij}) \right)$$



GNN structure

3. Compute distribution over variables



Outline (applied techniques)

1. GNNs overview
2. Integer programming with GNNs
 - i. Machine learning formulation
 - ii. Baselines**
 - iii. Experiments
 - iv. Additional research
3. Neural algorithmic alignment
4. Learning greedy heuristics with RL

Reliability pseudo-cost branching (RPB)

Rough idea:

- Goal: estimate $z(j) - z_i^+(j)$ w/o solving the LP with $x_i = 1$
- Estimate = avg change after setting $x_i = 1$ elsewhere in tree
This is the "pseudo-cost"
- "Reliability": do strong branching if estimate is "unreliable"
E.g., early in the tree

Default branching rule of SCIP (leading open-source solver):

$$\tilde{\Delta}_i^+(j) \cdot \tilde{\Delta}_i^-(j)$$

Estimate of $z(j) - z_i^+(j)$

Estimate of $z(j) - z_i^-(j)$

Technically,
 $\max\{\tilde{\Delta}_i^+(j), 10^{-6}\} \cdot \max\{\tilde{\Delta}_i^-(j), 10^{-6}\}$

Learning to rank approaches

Predict which variable **strong branching** would rank highest using models other than GNNs

- Khalil et al. [AAAI'16]:
Use learning-to-rank algorithm **SVM^{rank}** [Joachims, KDD'06]
- Hansknecht et al. [arXiv'18]:
Use learning-to-rank alg **lambdaMART** [Burgess, Learning'10]
- Alvarez et al. [INFORMS JoC'17]:
Use **regression trees**

Outline (applied techniques)

1. GNNs overview
2. Integer programming with GNNs
 - i. Machine learning formulation
 - ii. Baselines
 - iii. Experiments**
 - iv. Additional research
3. Neural algorithmic alignment
4. Learning greedy heuristics with RL

Set covering instances

Train and test on "easy" instances: 1000 columns, 500 rows

Model	Time	Wins	Nodes
Reliability	Runtime in seconds with a timeout of 1 hour	0/100	17 ± 13.7%
Regression trees	9.28 ± 4.9%	0/100	54 ± 20.8%
SVMrank	8.10 ± 3.8%	1/100	165 ± 8.2%
lambdaMART	7.19 ± 4.2%	14/100	167 ± 9.0%
GNN	6.59 ± 3.1%	85/100	134 ± 7.6%

Number instances with fastest runtime / number solved

Size of B&B tree

Set covering instances

Train and test on “easy” instances: 1000 columns, 500 rows

Model	Time	Wins	Nodes
Full strong branching	17.30±6.1%	0/100	17±13.7%
Reliability pseudo-cost	8.98±4.8%	0/100	54 ±20.8%
Regression trees	9.28±4.9%	0/100	187±9.4%
SVMrank	8.10±3.8%	1/100	165±8.2%
lambdaMART	7.19±4.2%	14/100	167±9.0%
GNN	6.59 ±3.1%	85 /100	134±7.6%

Set covering instances

GNN is **faster than SCIP** default VSP (reliability pseudo-cost)

Model	Time	Wins	Nodes
Full strong branching	17.30±6.1%	0/100	17±13.7%
Reliability pseudo-cost	8.98±4.8%	0/100	54 ±20.8%
Regression trees	9.28±4.9%	0/100	187±9.4%
SVMrank	8.10±3.8%	1/100	165±8.2%
lambdaMART	7.19±4.2%	14/100	167±9.0%
GNN	6.59 ±3.1%	85 /100	134±7.6%

Set covering instances

Train: "easy"; test: "**hard**" instances w/ 1000 columns, 2000 rows

Model	Time	Wins	Nodes
Full strong branching	Timed out	0/0	N/A
Reliability pseudo-cost	1677.98 \pm 3.0%	4/65	47299 \pm 4.9%
Regression trees	2869.21 \pm 3.2%	0/35	59013 \pm 9.3%
SVMrank	2389.92 \pm 2.3%	0/47	42120 \pm 5.4%
lambdaMART	2165.96 \pm 2.0%	0/54	45319 \pm 3.4%
GNN	1489.91\pm3.3%	66/70	29981\pm4.9%

Set covering instances

Performance generalizes to **larger instances**

Model	Time	Wins	Nodes
Full strong branching	Timed out	0/0	N/A
Reliability pseudo-cost	1677.98 \pm 3.0%	4/65	47299 \pm 4.9%
Regression trees	2869.21 \pm 3.2%	0/35	59013 \pm 9.3%
SVMrank	2389.92 \pm 2.3%	0/47	42120 \pm 5.4%
lambdaMART	2165.96 \pm 2.0%	0/54	45319 \pm 3.4%
GNN	1489.91\pm3.3%	66/70	29981\pm4.9%

Set covering instances

Similar results for auction design & facility location problems

Outline (applied techniques)

1. GNNs overview
2. Integer programming with GNNs
 - i. Machine learning formulation
 - ii. Baselines
 - iii. Experiments
 - iv. Additional research**
3. Neural algorithmic alignment
4. Learning greedy heuristics with RL

Additional research

CPU-friendly approaches

Gupta et al., NeurIPS'20

Bipartite representation inspired many follow-ups

Nair et al., '20; Sonnerat et al., '21; Wu et al., NeurIPS'21; Huang et al. ICML'23; ...

Survey on *Combinatorial Optimization & Reasoning w/ GNNs*:

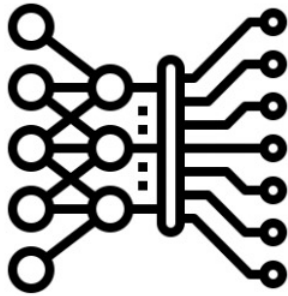
Cappart, Chételat, Khalil, Lodi, Morris, Veličković, JMLR'23

Outline (applied techniques)

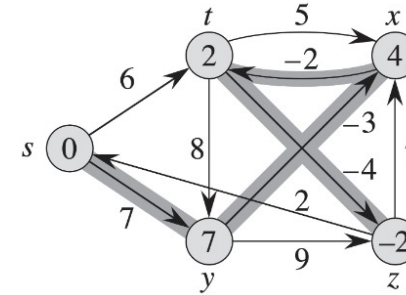
1. GNNs overview
2. Integer programming with GNNs
- 3. Neural algorithmic alignment**
4. Learning greedy heuristics with RL

Veličković, Ying, Padovano, Hadsell, Blundell, ICLR'20
Cappart, Chételat, Khalil, Lodi, Morris, Veličković, JMLR'23

Problem-solving approaches



- + Operate on raw inputs
- + Generalize on noisy conditions
- + Models reusable across tasks
- Require big data
- Unreliable when extrapolating
- Lack of interpretability



- + Trivially strong generalization
- + Compositional (subroutines)
- + Guaranteed correctness
- + Interpretable operations
- Input must match spec
- Not robust to task variations

Is it possible to get the best of both worlds?

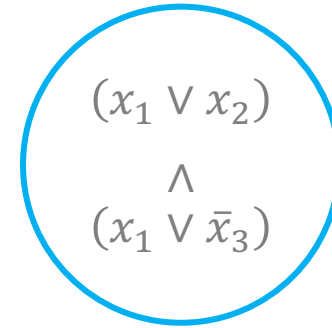
GNNs + combinatorial optimization

Lots of awesome research! E.g.,



Traveling salesman problem

E.g., Vinyals et al., '15; Joshi et al., '19; ...



Boolean satisfiability

E.g., Selsam et al., '19; Cameron et al., '20; ...

This section: Neural graph algorithm execution

Aligns well with theoretical sections of this tutorial

Neural graph algorithm execution

Key observation: Many algorithms share related **subroutines**
E.g. Bellman-Ford & BFS enumerate sets of edges adjacent to a node

Neural graph algorithm execution

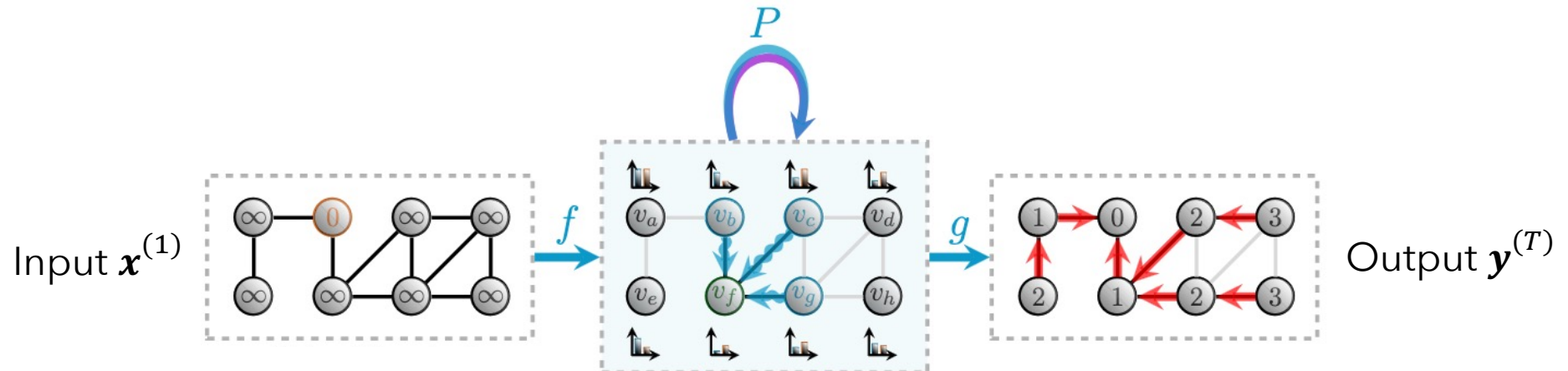
💡 Learn several algorithms **simultaneously**

If we already have a classical algorithm for the problem...

Why not just run that algorithm?

Will answer soon, but first: a few words on the pipeline

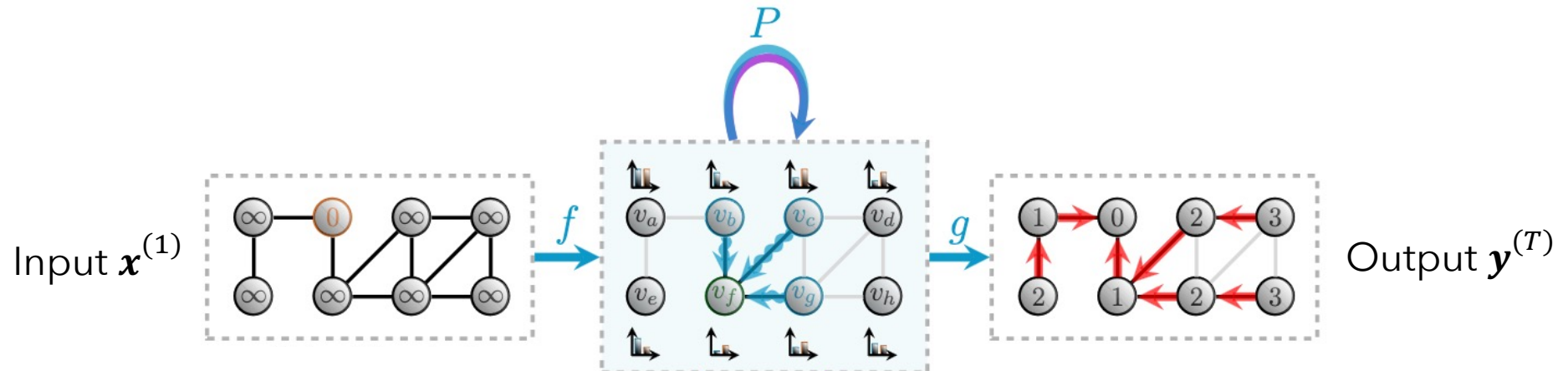
Neural algorithmic pipeline



Encoder network f

- E.g., makes sure input is in correct dimension for next step

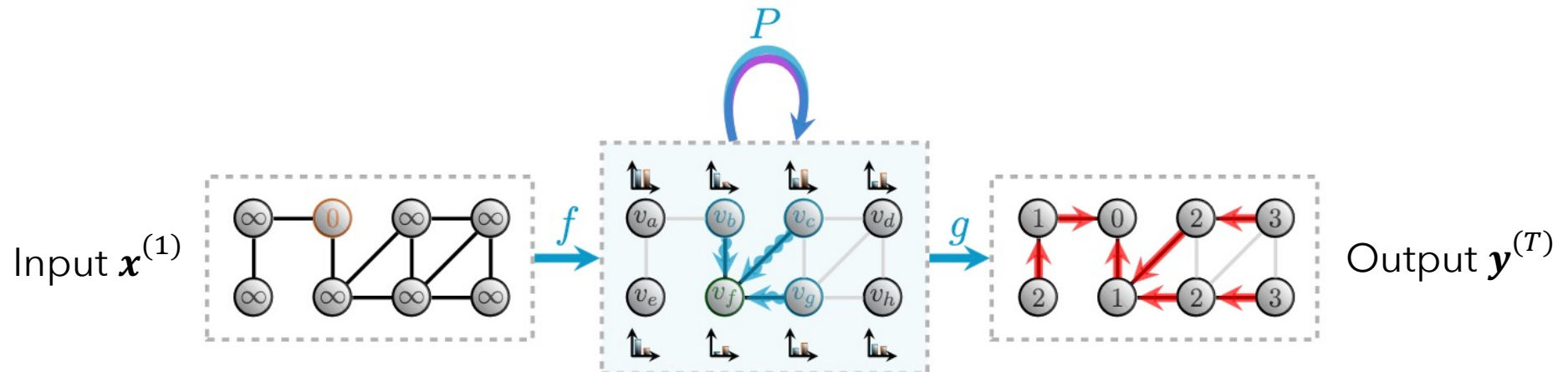
Neural algorithmic pipeline



Processor network P

- Graph neural network
- Run multiple times (termination determined by a NN)

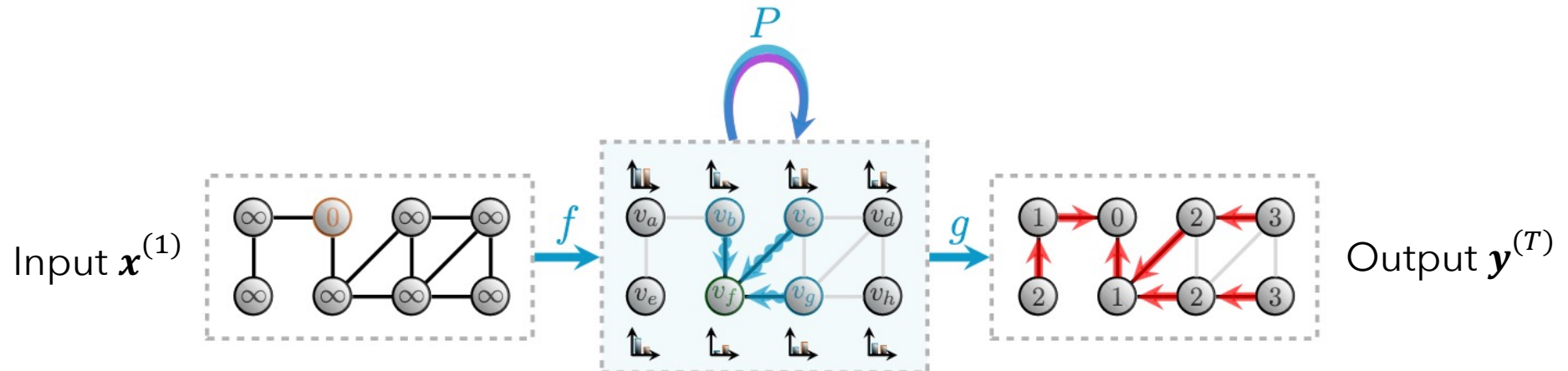
Neural algorithmic pipeline



Decoder network g

- Transform's GNNs output into algorithmic output

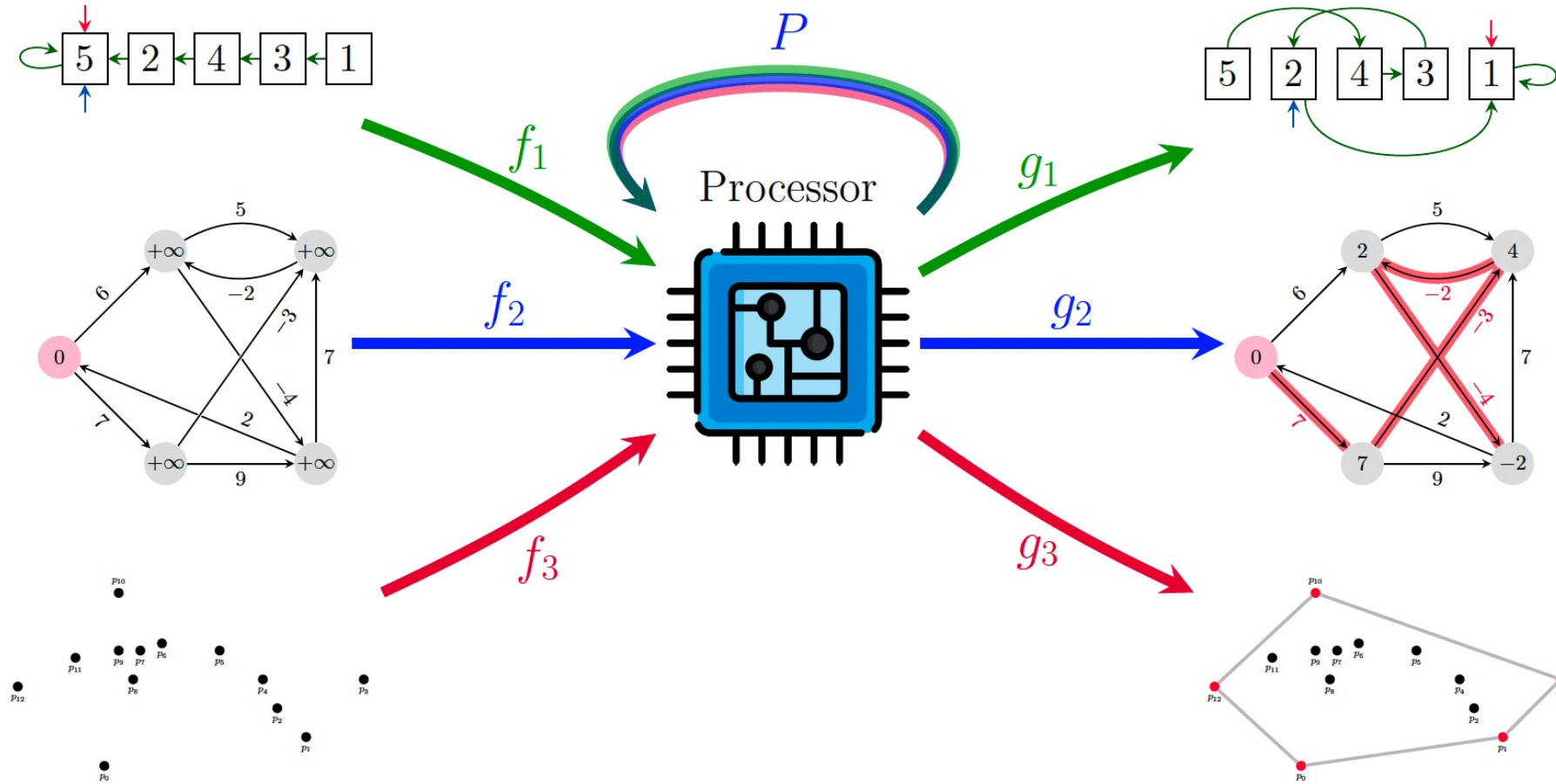
Neural algorithmic pipeline



Multi-task approach

- Learn a **single** processor network P for related problems
- Learn **task-specific** encoder, decoder functions f_A, g_A

Neural algorithmic pipeline



Why use GNNs for algorithm design?

If we're just teaching a NN to **imitate** a classical algorithm...

Why not just run that algorithm?

Why use GNNs for algorithm design?

Classical algorithms are designed with **abstraction** in mind
Enforce their inputs to conform to stringent preconditions

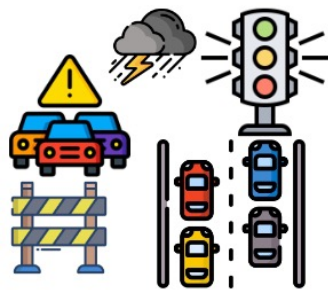
However, we design algorithms to solve **real-world** problems!



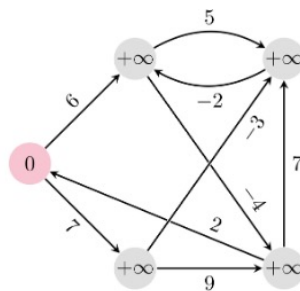
Natural inputs

Why use GNNs for algorithm design?

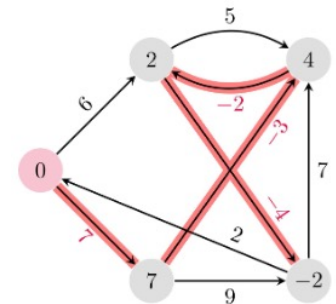
- Assume we have real-world inputs
...but algorithm only admits abstract inputs
- Could try **manually** converting from one input to another



Natural inputs



Abstract inputs



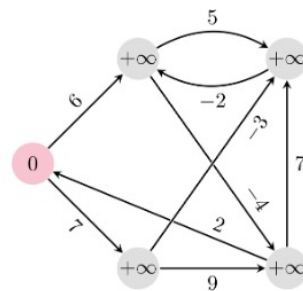
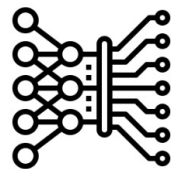
Abstract outputs

Why use GNNs for algorithm design?

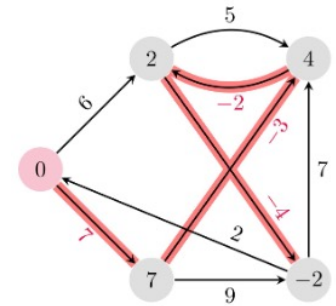
- Alternatively, **replace** human feature extractor with NN
 - Still apply same combinatorial algorithm
- Issue: algorithms typically perform **discrete optimization**
 - Doesn't play nicely with **gradient-based** optimization of NNs



Natural inputs



Abstract inputs



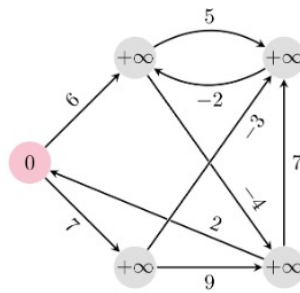
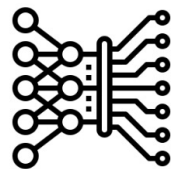
Abstract outputs

Why use GNNs for algorithm design?

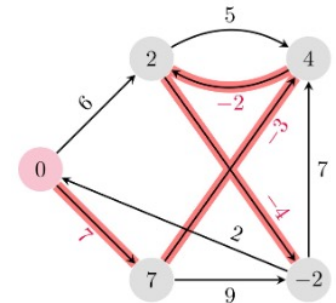
- Second (more fundamental) issue: **data efficiency**
 - Real-world data is often incredibly rich
 - We still have to compress it down to scalar values
- The algorithmic solver commits to using this scalar
Assumes it is perfect!



Natural inputs



Abstract inputs



Abstract outputs

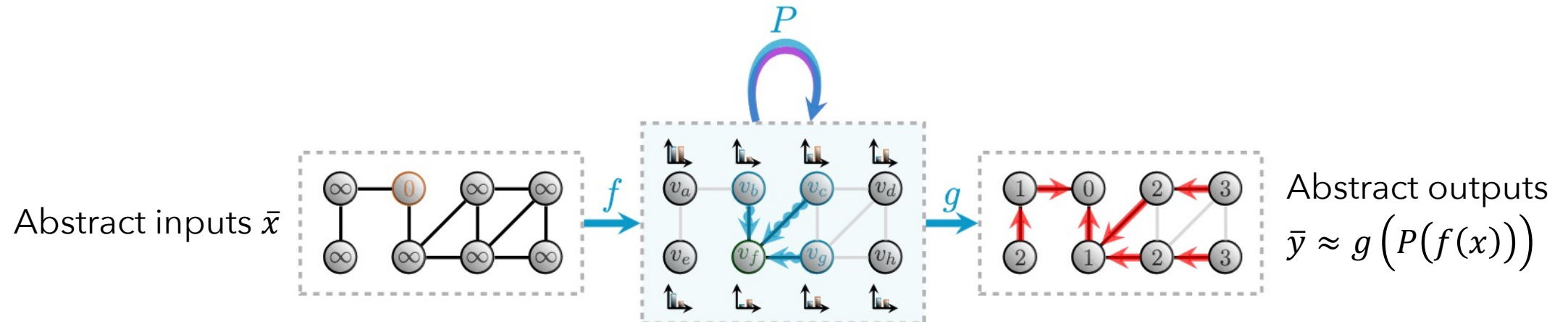
Why use GNNs for algorithm design?

- Second (more fundamental) issue: **data efficiency**
 - Real-world data is often incredibly rich
 - We still have to compress it down to scalar values
- The algorithmic solver commits to using this scalar
Assumes it is perfect!

If there's insufficient training data to estimate the scalars:

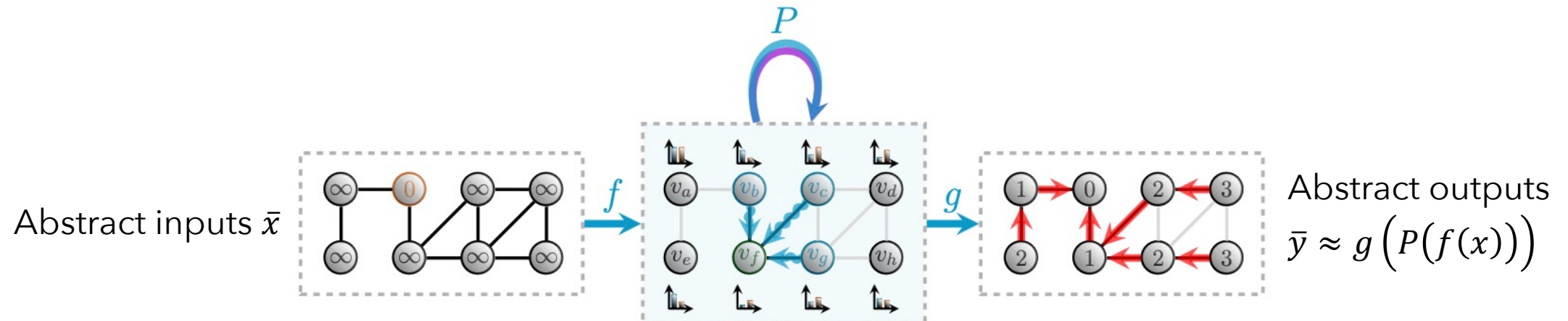
- Alg will give a **perfect solution**
- ...but in a **suboptimal environment**

Neural algorithmic pipeline



1. On abstract inputs, learn encode-process-decode functions

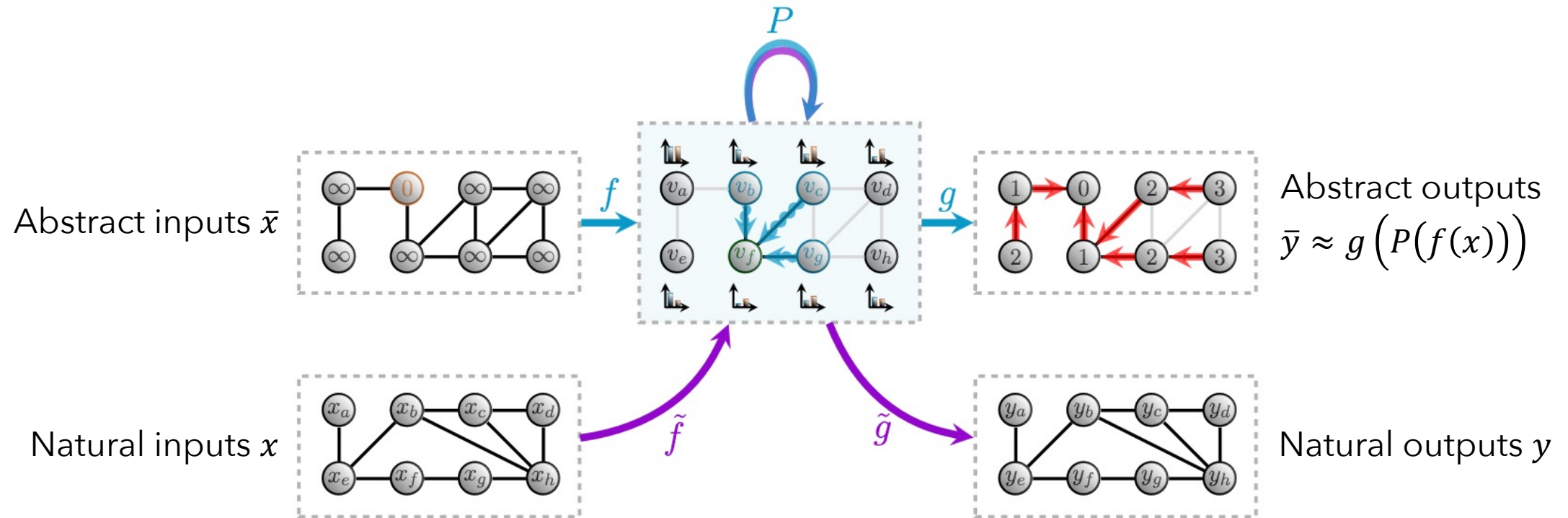
Neural algorithmic pipeline



After training on abstract inputs, processor P :

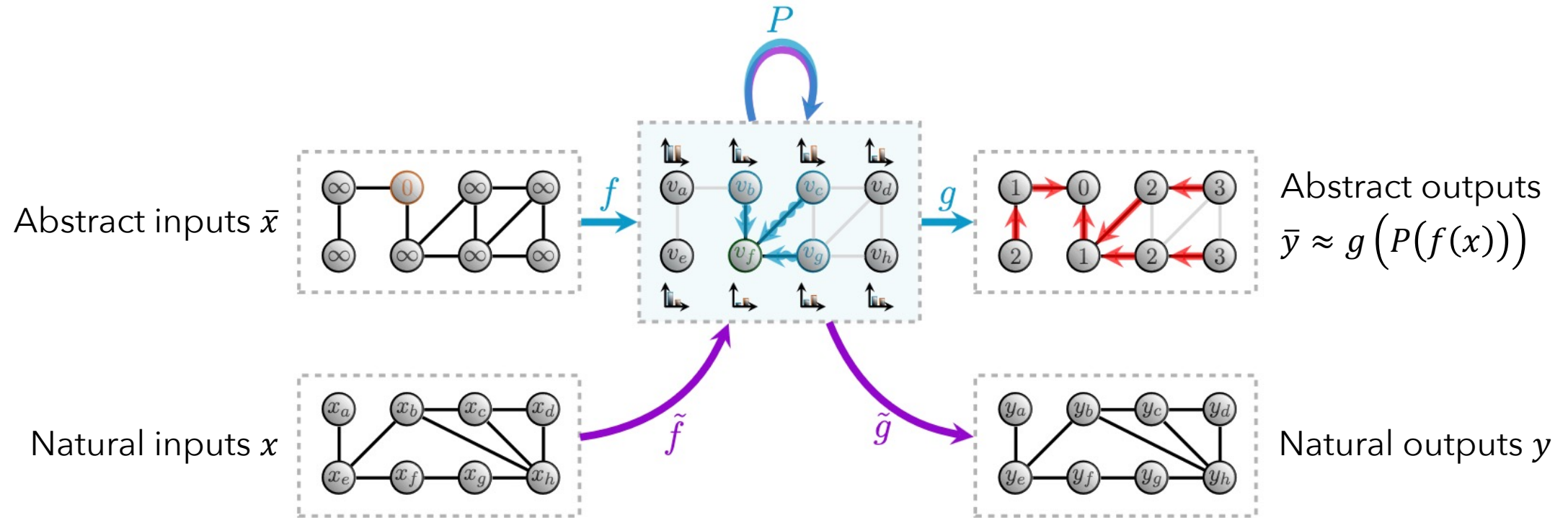
1. Admits useful gradients
2. Operates over high-dim latent space (better use of data)

Neural algorithmic pipeline



2. Set up encode-decode functions for natural inputs/outputs

Neural algorithmic pipeline



3. Learn parameters using loss that compares $\tilde{g} \left(P \left(\tilde{f}(x) \right) \right)$ to y

Outline (applied techniques)

1. GNNs overview
2. Integer programming with GNNs
3. Neural algorithmic alignment
 - i. Motivation
 - ii. Example algorithms**
 - iii. Experiments
 - iv. Understanding max-aggregation
 - v. Additional research
4. Learning greedy heuristics with RL

Breadth-first search

- Source node s

- Initial input $x_i^{(1)} = \begin{cases} 1 & \text{if } i = s \\ 0 & \text{if } i \neq s \end{cases}$

- Node is reachable from s if any of its neighbors are reachable:

$$x_i^{(t+1)} = \begin{cases} 1 & \text{if } x_i^{(t)} = 1 \\ 1 & \text{if } \exists j \text{ s.t. } (j, i) \in E \text{ and } x_j^{(t)} = 1 \\ 0 & \text{else} \end{cases}$$

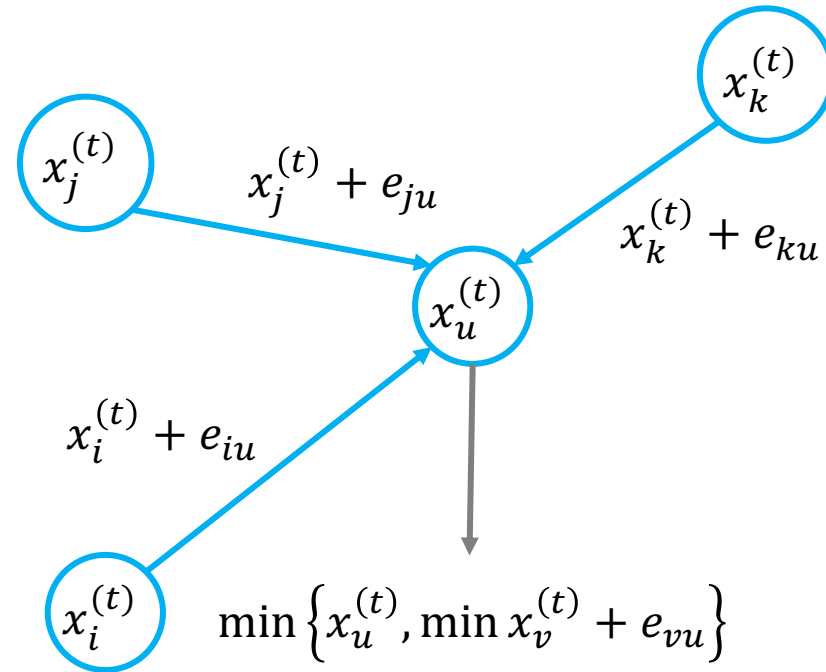
- Algorithm output at round t : $y_i^{(t)} = x_i^{(t+1)}$

Bellman-Ford (shortest path)

- Source node s
- Initial input $x_i^{(1)} = \begin{cases} 0 & \text{if } i = s \\ \infty & \text{if } i \neq s \end{cases}$
- Node is reachable from s if any of its neighbors are reachable
Update distance to node as minimal way to reach neighbors

$$x_i^{(t+1)} = \min \left\{ x_i^{(t)}, \min_{(j,i) \in E} x_j^{(t)} + e_{ji} \right\}$$

Bellman-Ford: Message passing

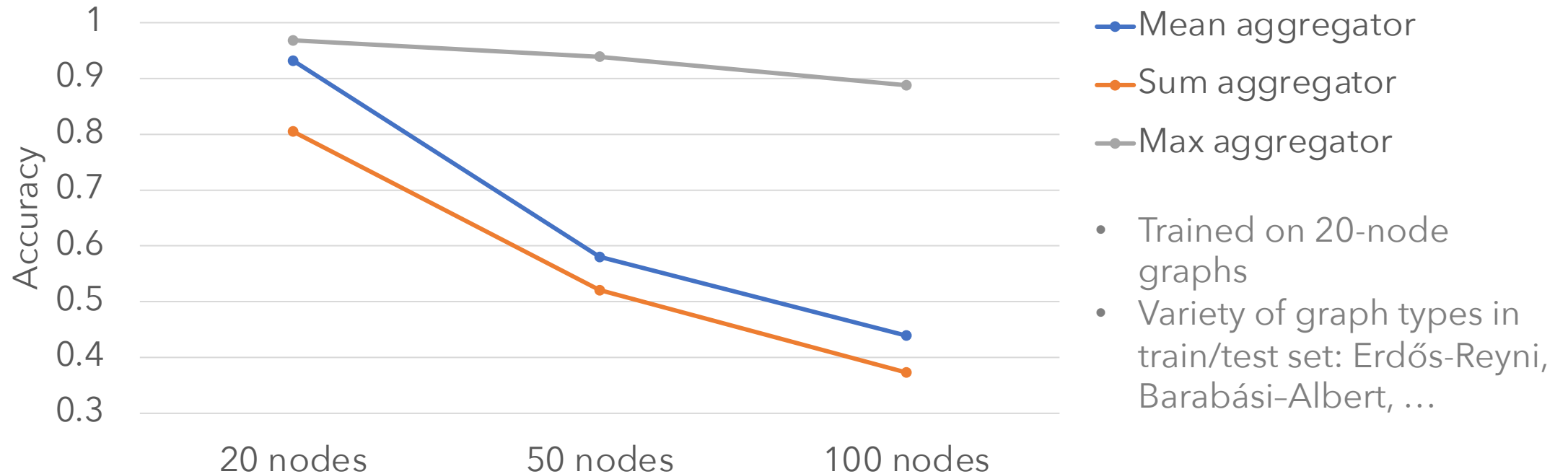


Key idea (roughly speaking): Train GNN so that $\mathbf{h}_u^{(t)} \approx x_u^{(t)}, \forall t$
(Really, so that a function of $\mathbf{h}_u^{(t)} \approx x_u^{(t)}$)

Outline (applied techniques)

1. GNNs overview
2. Integer programming with GNNs
3. Neural algorithmic alignment
 - i. Motivation
 - ii. Example algorithms
 - iii. Experiments**
 - iv. Understanding max-aggregation
 - v. Additional research
4. Learning greedy heuristics with RL

Shortest-path predecessor prediction



Improvement of max-aggregator increases with size

It **aligns** better with underlying algorithm [Xu et al., ICLR'20]

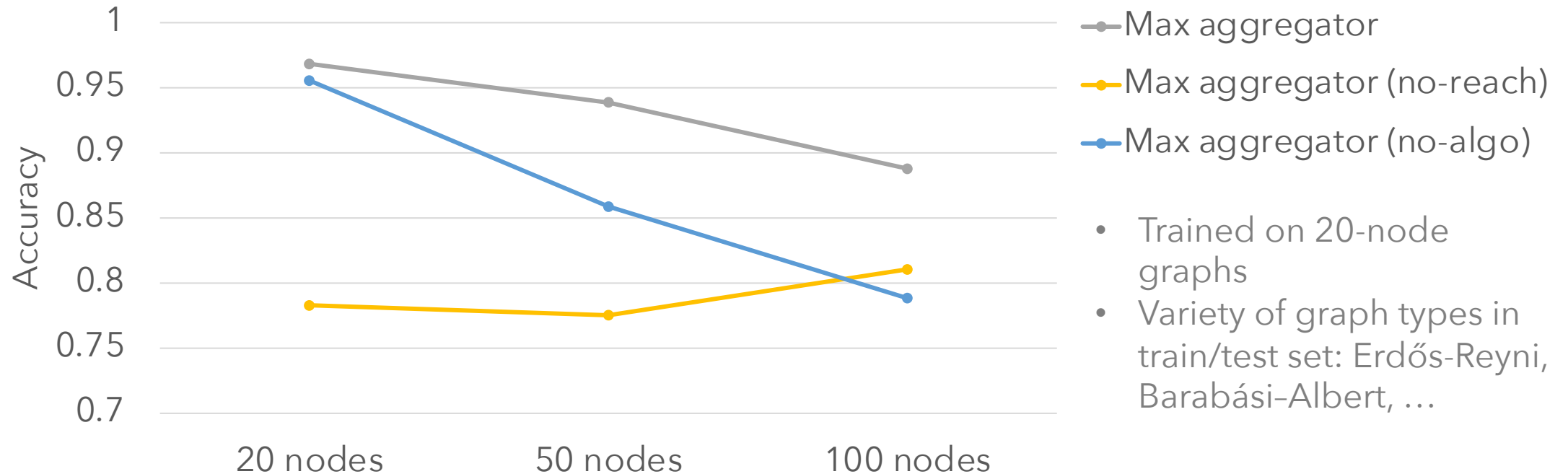
Learning multiple algorithms

Learn to execute both BFS and Bellman-Ford **simultaneously**

Comparisons

- (*no-reach*): Learn Bellman-Ford alone
 - Doesn't simultaneously learn reachability
- (*no-algo*):
 - Don't supervise intermediate steps
 - Learn predecessors directly from input $x_i^{(1)}$

Shortest-path predecessor prediction



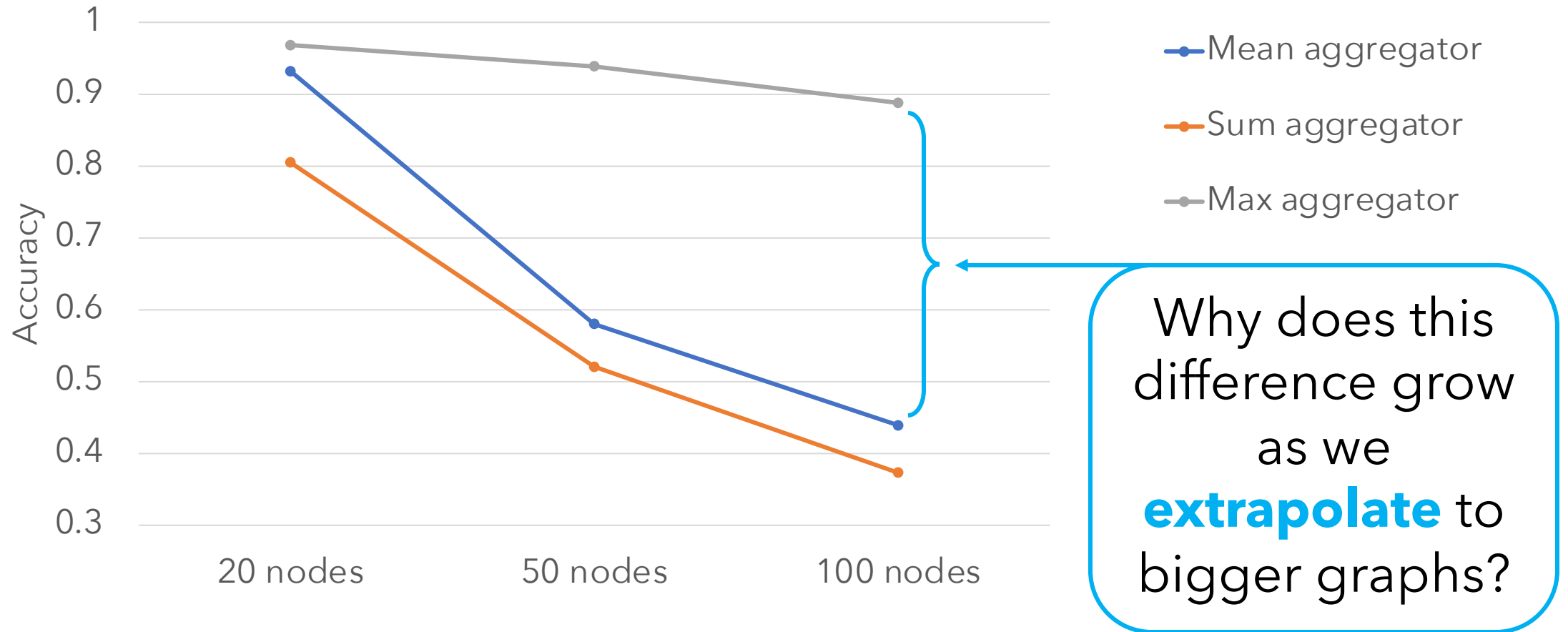
- **(no-reach) results:** positive knowledge transfer
- **(no-algo) results:** benefit of supervising intermediate steps

Outline (applied techniques)

1. GNNs overview
2. Integer programming with GNNs
3. Neural algorithmic alignment
 - i. Motivation
 - ii. Example algorithms
 - iii. Experiments
 - iv. Understanding max-aggregation**
 - v. Additional research
4. Learning greedy heuristics with RL

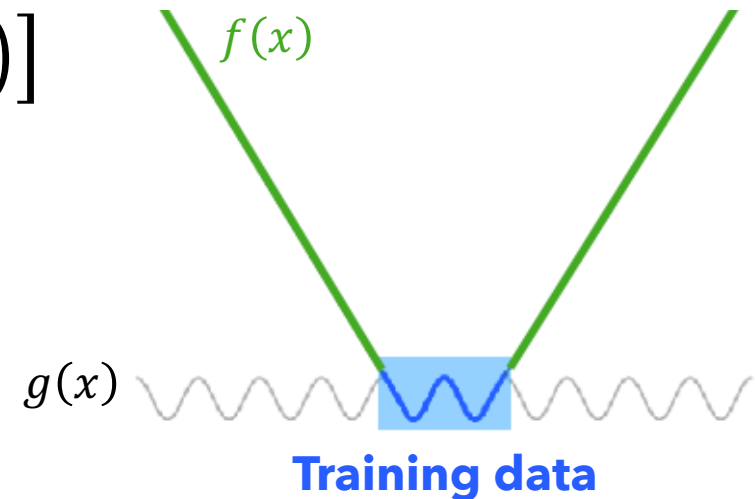
Xu, Zhang, Li, Du, Kawarabayashi, Jegelka, ICLR'21

Shortest-path predecessor prediction

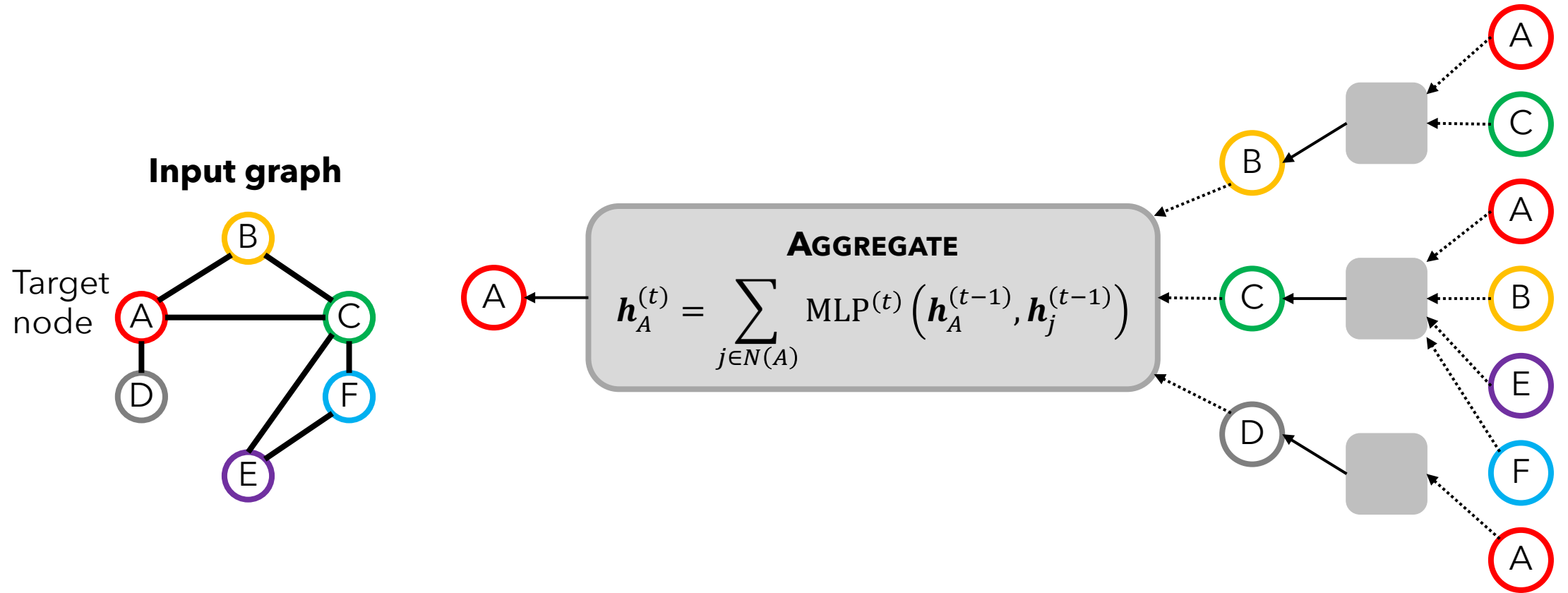


Extrapolation error

- $f: \mathcal{X} \rightarrow \mathbb{R}$ is a model trained on $\{(x_i, y_i)\}_{i=1}^n \subset \mathcal{D}$
 $y_i = g(x_i)$ for some ground-truth function g
- \mathcal{P} is a **distribution over $\mathcal{X} \setminus \mathcal{D}$**
- $\ell: \mathbb{R} \times \mathbb{R} \rightarrow \mathbb{R}$ is a loss function
- **Extrapolation error:** $\mathbb{E}_{x \sim \mathcal{P}} [\ell(f(x), g(x))]$



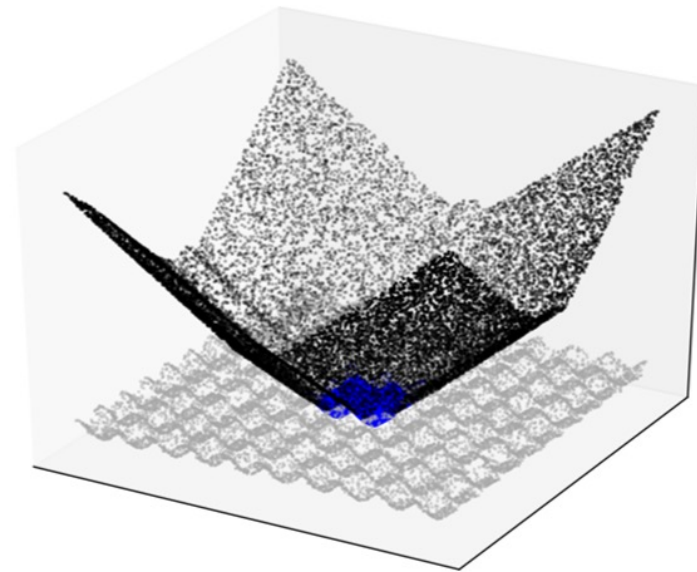
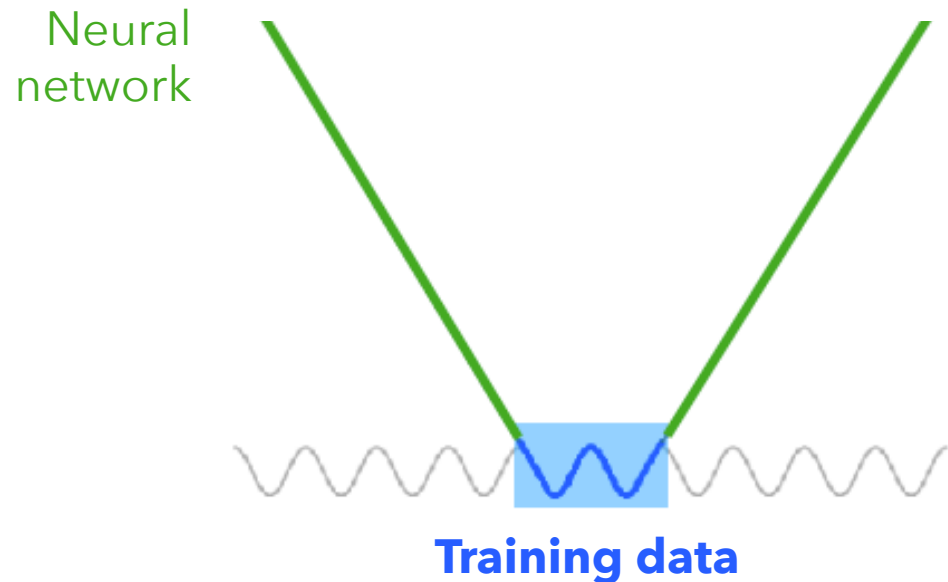
Aggregation functions



ReLU MLP extrapolate linearly

Theorem [Xu et al., ICLR'21, informal]:

- $f: \mathbb{R}^d \rightarrow \mathbb{R}$, a 2-layer ReLU MLP trained w/ gradient descent
- Along any direction $\mathbf{v} \in \mathbb{R}^d$, f approaches a **linear** function



ReLU MLP extrapolate linearly

Theorem [Xu et al., ICLR'21, informal]:

- $f: \mathbb{R}^d \rightarrow \mathbb{R}$, a 2-layer ReLU MLP trained w/ gradient descent
- Along any direction $\mathbf{v} \in \mathbb{R}^d$, f approaches a **linear** function
- More formally, let $\mathbf{x} = t\mathbf{v}$
 - Then $f(\mathbf{x} + h\mathbf{v}) - f(\mathbf{x}) = f(t\mathbf{v} + h\mathbf{v}) - f(t\mathbf{v}) \rightarrow \beta_{\mathbf{v}}h$
at a rate $O\left(\frac{1}{t}\right)$

Implications for GNNs

Shortest path: $x_i^{(t)} = \min \left\{ x_i^{(t-1)}, \min_{(j,i) \in E} x_j^{(t-1)} + e_{ji} \right\}$

GNN: $\mathbf{h}_i^{(t)} = \sum_{j \in N(i)} \text{MLP} \left(\mathbf{h}_i^{(t-1)}, \mathbf{h}_j^{(t-1)} \right)$

MLP must learn a **non-linearity**

Implications for GNNs

Shortest path:
$$x_i^{(t)} = \min \left\{ x_i^{(t-1)}, \min_{(j,i) \in E} x_j^{(t-1)} + e_{ji} \right\}$$

GNN:

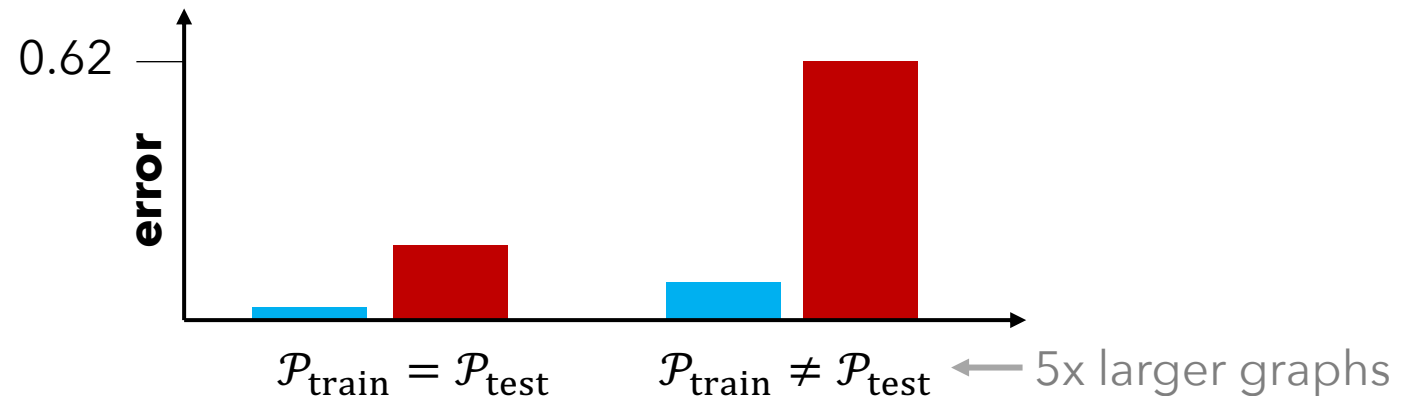
$$\mathbf{h}_i^{(t)} = \sum_{j \in N(i)} \text{MLP} \left(\mathbf{h}_i^{(t-1)}, \mathbf{h}_j^{(t-1)} \right)$$

GNN 2:

$$\mathbf{h}_i^{(t)} = \max_{j \in N(i)} \text{MLP} \left(\mathbf{h}_i^{(t-1)}, \mathbf{h}_j^{(t-1)} \right)$$

Predicting shortest path predecessor:

[Veličković et al. ICLR'20]



Outline (applied techniques)

1. GNNs overview
2. Integer programming with GNNs
3. Neural algorithmic alignment
 - i. Motivation
 - ii. Example algorithms
 - iii. Experiments
 - iv. Understanding max-aggregation
 - v. Additional research**
4. Learning greedy heuristics with RL

Additional research

Lots of research in the past few years! E.g.:

- How to achieve **algorithmic alignment** & **theory guarantees**
 - Xu et al., ICLR'20; Dudzik, Veličković, NeurIPS'22
- **CLRS** benchmark
 - Sorting, searching, dynamic programming, graph algorithms, etc.
 - Veličković et al. ICML'22; Ibarz et al. LoG'22; Bevilacqua et al. ICML'23
- **Primal-dual** algorithms
 - Numeroso et al., ICLR'23

Outline (applied techniques)

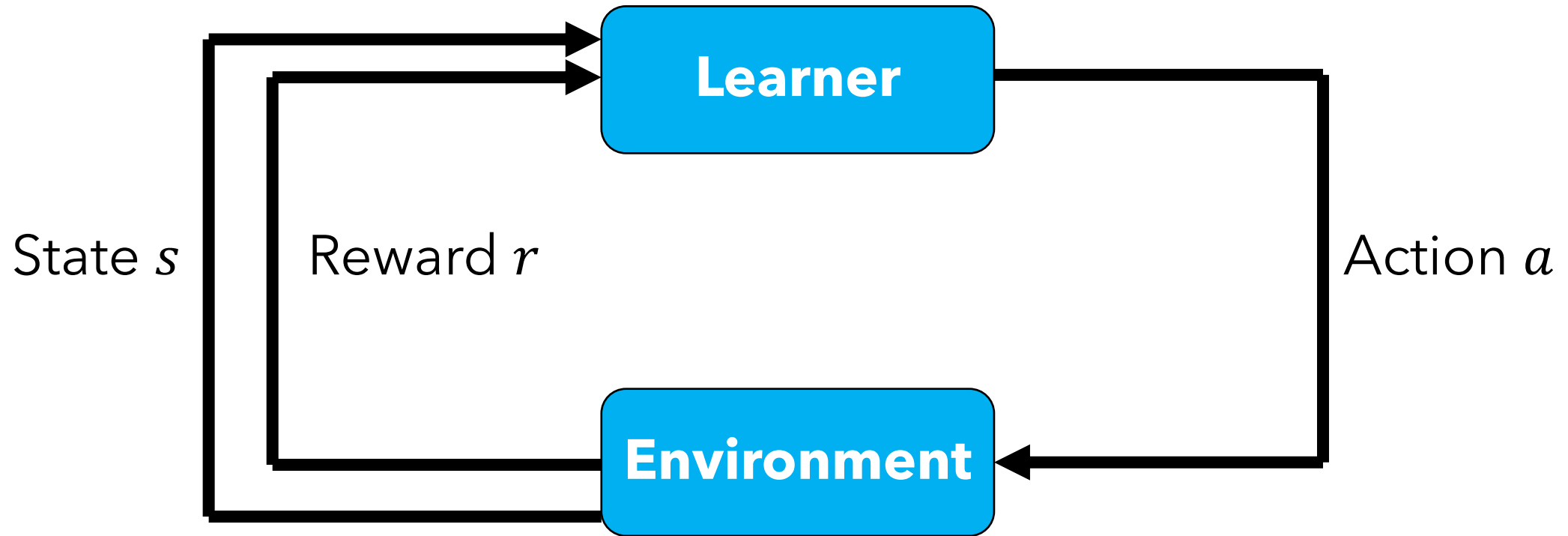
1. GNNs overview
2. Integer programming with GNNs
3. Neural algorithmic alignment
- 4. Learning greedy heuristics with RL**

Dai, Khalil, Zhang, Dilkina, Song; NeurIPS'17

Outline (applied techniques)

1. GNNs overview
2. Integer programming with GNNs
3. Neural algorithmic alignment
4. Learning greedy heuristics with RL
 - i. Reinforcement learning refresher**
 - a. Markov decision processes**
 - b. Reinforcement learning
 - ii. Overview: RL for combinatorial optimization
 - iii. Examples: Min vertex cover and max cut
 - iv. RL formulation
 - v. Experiments

Learner interaction with environment



Markov decision processes

S : set of states (assumed for now to be discrete)

A : set of actions

Transition probability distribution $P(s' | s, a)$

Probability of entering state s' from state s after taking action a

Reward function $R: S \rightarrow \mathbb{R}$

Goal: Policy $\pi: S \rightarrow A$ that maximizes total (discounted) reward

Policies and value functions

Value function for a policy:

Expected sum of discounted rewards

$$V^\pi(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \mid s_0 = s, a_t = \pi(s_t), (s_{t+1} | s_t, a_t) \sim P \right]$$

Discount factor

$$= R(s) + \gamma \sum_{s' \in \mathcal{S}} P(s' | s, \pi(s)) V^\pi(s') \quad \text{(Bellman equation)}$$

Optimal policy and value function

Optimal policy π^* achieves the highest value for every state

$$V^{\pi^*}(s) = \max_{\pi} V^{\pi}(s)$$

Several different ways to find π^*

- Value iteration
- Policy iteration

Outline (applied techniques)

1. GNNs overview
2. Integer programming with GNNs
3. Neural algorithmic alignment
4. Learning greedy heuristics with RL
 - i. Reinforcement learning refresher
 - a. Markov decision processes
 - b. Reinforcement learning**
 - ii. Overview: RL for combinatorial optimization
 - iii. Examples: Min vertex cover and max cut
 - iv. RL formulation
 - v. Experiments

Challenge of RL

MDP (S, A, P, R):

- S : set of states (assumed for now to be discrete)
- A : set of actions
- Transition probability distribution $P(s_{t+1} \mid s_t, a_t)$
- Reward function $R: S \rightarrow \mathbb{R}$

RL twist: We don't know P or R , or too big to enumerate

Q-learning

Q functions:

Like value functions but defined over state-action pairs

$$Q^\pi(s, a) = R(s) + \gamma \sum_{s' \in S} P(s' | s, a) Q^\pi(s', \pi(s'))$$

I.e., Q function is the value of:

1. Starting in state s
2. Taking action a
3. Then acting according to π

Q-learning

Q function of the optimal policy π^* :

$$\begin{aligned} Q^*(s, a) &= R(s) + \gamma \sum_{s' \in S} P(s' | s, a) \max_{a'} Q^*(s', a') \\ &= R(s) + \gamma \sum_{s' \in S} P(s' | s, a) V^{\pi^*}(s') \end{aligned}$$

Q^* is the value of:

1. Starting in state s
2. Taking action a
3. Then acting optimally

Q-learning

(High-level) Q-learning algorithm

initialize $\hat{Q}(s, a) \leftarrow 0, \forall s, a$

repeat

Observe current state s and reward r

Take action $a = \operatorname{argmax} \hat{Q}(s, \cdot)$ and observe next state s'

Improve estimate \hat{Q} based on s, r, a, s'

Can use *function approximation* to represent \hat{Q} compactly

$$\hat{Q}(s, a) = f_{\theta}(s, a)$$

Outline (applied techniques)

1. GNNs overview
2. Integer programming with GNNs
3. Neural algorithmic alignment
4. Learning greedy heuristics with RL
 - i. Reinforcement learning refresher
 - ii. Overview: RL for combinatorial optimization**
 - iii. Examples: Min vertex cover and max cut
 - iv. RL formulation
 - v. Experiments

RL for combinatorial optimization

Tons of research in this area

Travelling salesman

Bello et al., ICLR'17; Dai et al., NeurIPS'17;
Nazari et al., NeurIPS'18; ...

Bin packing

Hu et al., '17; Laterre et al., '18; Cai et al.,
DRL4KDD'19; Li et al., '20; ...

Maximum cut

Dai et al., NeurIPS'17; Cappart et al.,
AAAI'19; Barrett et al., AAAI'20; ...

Minimum vertex cover

Dai et al., NeurIPS'17; Song et al., UAI'19; ...

This section: Example of a pioneering work in this space

Overview

Goal: use RL to learn new *greedy strategies* for graph problems
Feasible solution constructed by successively adding nodes to solution

Input: Graph $G = (V, E)$, weights $w(u, v)$ for $(u, v) \in E$

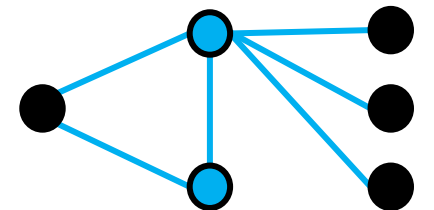
RL state representation: Graph embedding

Outline (applied techniques)

1. GNNs overview
2. Integer programming with GNNs
3. Neural algorithmic alignment
4. Learning greedy heuristics with RL
 - i. Reinforcement learning refresher
 - ii. Overview: RL for combinatorial optimization
 - iii. Examples: Min vertex cover and max cut**
 - iv. RL formulation
 - v. Experiments

Minimum vertex cover

Find smallest vertex subset such that each edge is covered

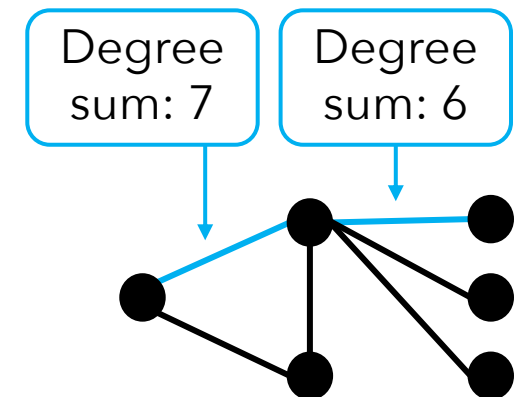


Minimum vertex cover

Find smallest vertex subset such that each edge is covered

2-approximation:

Greedily add vertices of edge with **maximum degree sum**



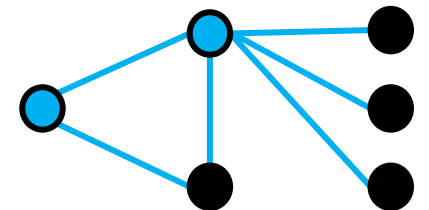
Minimum vertex cover

Find smallest vertex subset such that each edge is covered

2-approximation:

Greedily add vertices of edge with maximum degree sum

Scoring function that guides greedy algorithm



Maximum cut

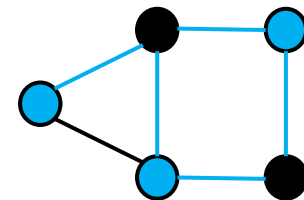
Find partition $(S, V \setminus S)$ of nodes that maximizes

$$\sum_{(u,v) \in C} w(u,v)$$

where $C = \{(u,v) \in E : u \in S, v \notin S\}$

If $w(u,v) = 1$ for all $(u,v) \in E$:

$$\sum_{(u,v) \in C} w(u,v) = 5$$



Maximum cut

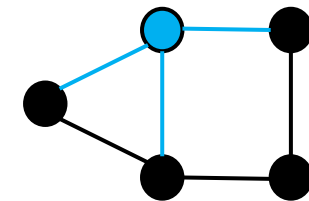
Find partition $(S, V \setminus S)$ of nodes that maximizes

$$\sum_{(u,v) \in C} w(u,v)$$

where $C = \{(u, v) \in E : u \in S, v \notin S\}$

Greedy: move node from one side of cut to the other

Move node that results in the largest improvement in cut weight



Maximum cut

Find partition $(S, V \setminus S)$ of nodes that maximizes

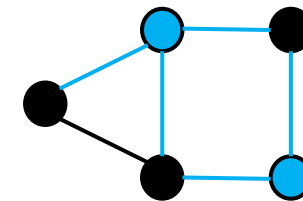
$$\sum_{(u,v) \in C} w(u,v)$$

where $C = \{(u, v) \in E : u \in S, v \notin S\}$

Greedy: move node from one side of cut to the other

Move node that results in the largest improvement in cut weight

Scoring function that guides greedy algorithm



Outline (applied techniques)

1. GNNs overview
2. Integer programming with GNNs
3. Neural algorithmic alignment
4. Learning greedy heuristics with RL
 - i. Reinforcement learning refresher
 - ii. Overview: RL for combinatorial optimization
 - iii. Examples: Min vertex cover and max cut
 - iv. RL formulation**
 - v. Experiments

RL for combinatorial optimization

Goal: learn a scoring function to guide greedy algorithm

Problem

Greedy operation

Min vertex cover

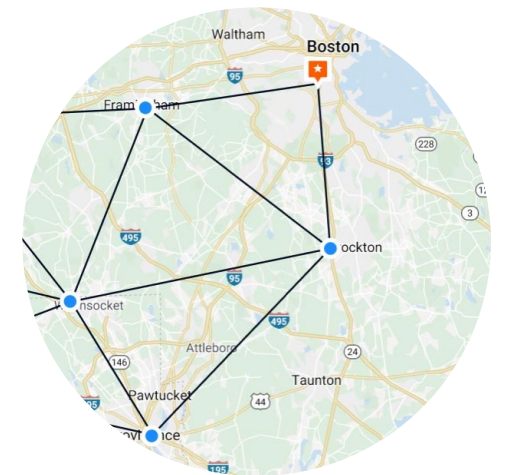
Insert node into cover

Max cut

Insert node into subset

Traveling salesman

Insert node into sub-tour



RL for combinatorial optimization

Greedy algorithm **Reinforcement learning**

Partial solution

State

Scoring function

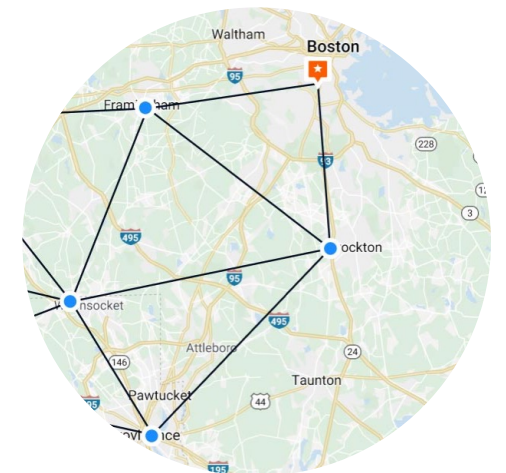
Q-function

Select best node

Greedy policy

Repeat until all edges are covered:

1. Compute node scores
2. Select best node with respect to score
3. Add best node to partial solution



Reinforcement learning formulation

State:

- *Goal*: encode partial solution $S = (v_1, v_2, \dots, v_{|S|})$, $v_i \in V$

E.g., nodes in independent set, nodes on one side of cut

Reinforcement learning formulation

State:

- *Goal*: encode partial solution $S = (v_1, v_2, \dots, v_{|S|}), v_i \in V$
- Use GNN to compute graph embedding μ

$$\text{Initial node features } x_v = \begin{cases} 1 & \text{if } v \in S \\ 0 & \text{else} \end{cases}$$

Action: Choose vertex $v \in V \setminus S$ to add to solution

Transition (deterministic): For chosen $v \in V \setminus S$, set $x_v = 1$

Reinforcement learning formulation

Reward: $r(S, v)$ is change in objective when transition $S \rightarrow (S, v)$

Policy (deterministic): $\pi(v|S) = \begin{cases} 1 & \text{if } v = \operatorname{argmax}_{v' \in S} \hat{Q}(\mu, v') \\ 0 & \text{else} \end{cases}$

Outline (applied techniques)

1. GNNs overview
2. Integer programming with GNNs
3. Neural algorithmic alignment
4. Learning greedy heuristics with RL
 - i. Reinforcement learning refresher
 - ii. Overview: RL for combinatorial optimization
 - iii. Examples: Min vertex cover and max cut
 - iv. RL formulation
 - v. Experiments**

Min vertex cover

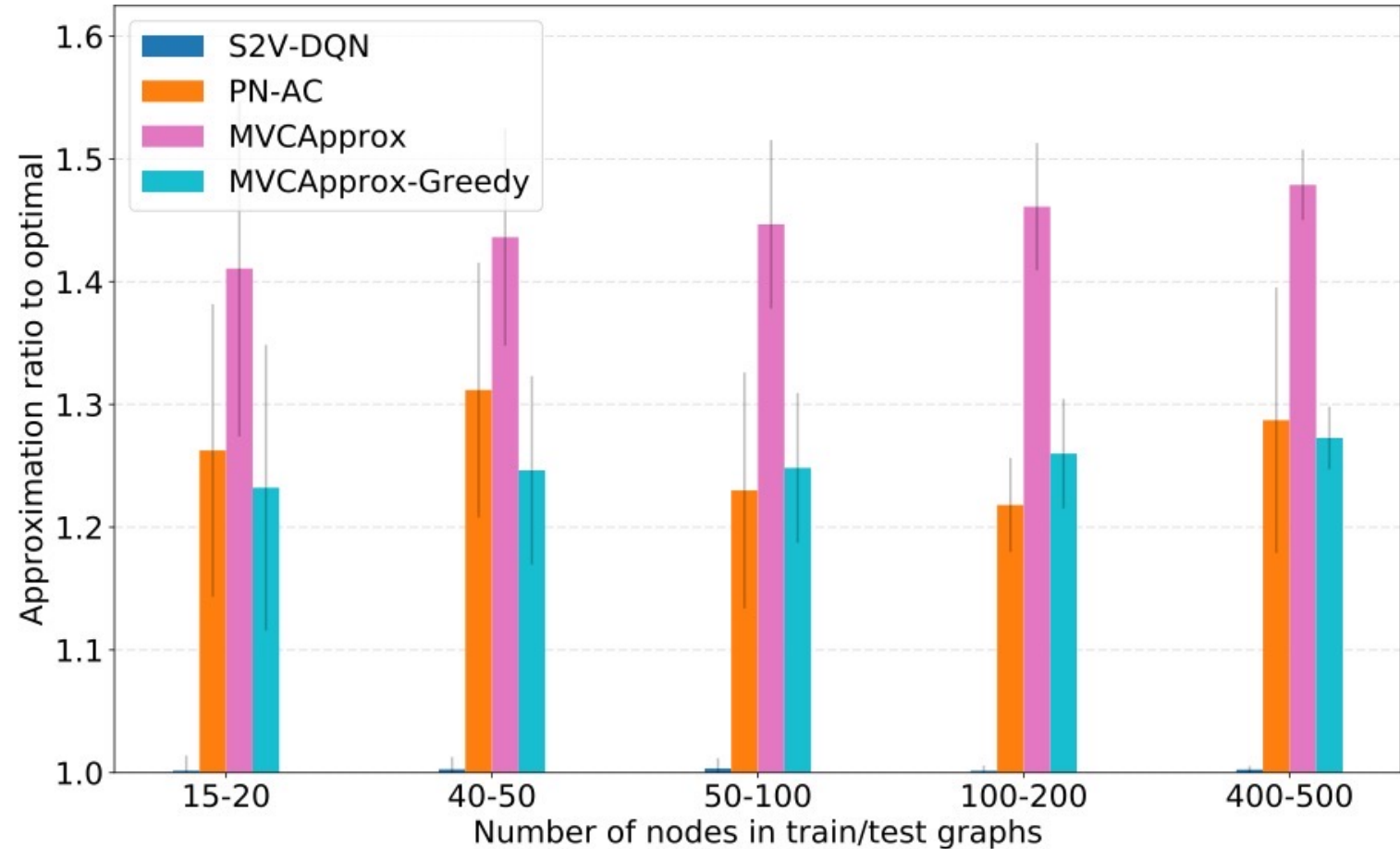
Barabasi-Albert
random graphs

Paper's approach

Another DL approach
[Bello et al., arXiv'16]

2-approximation
algorithm

Greedy algorithm
from first few slides



Max cut

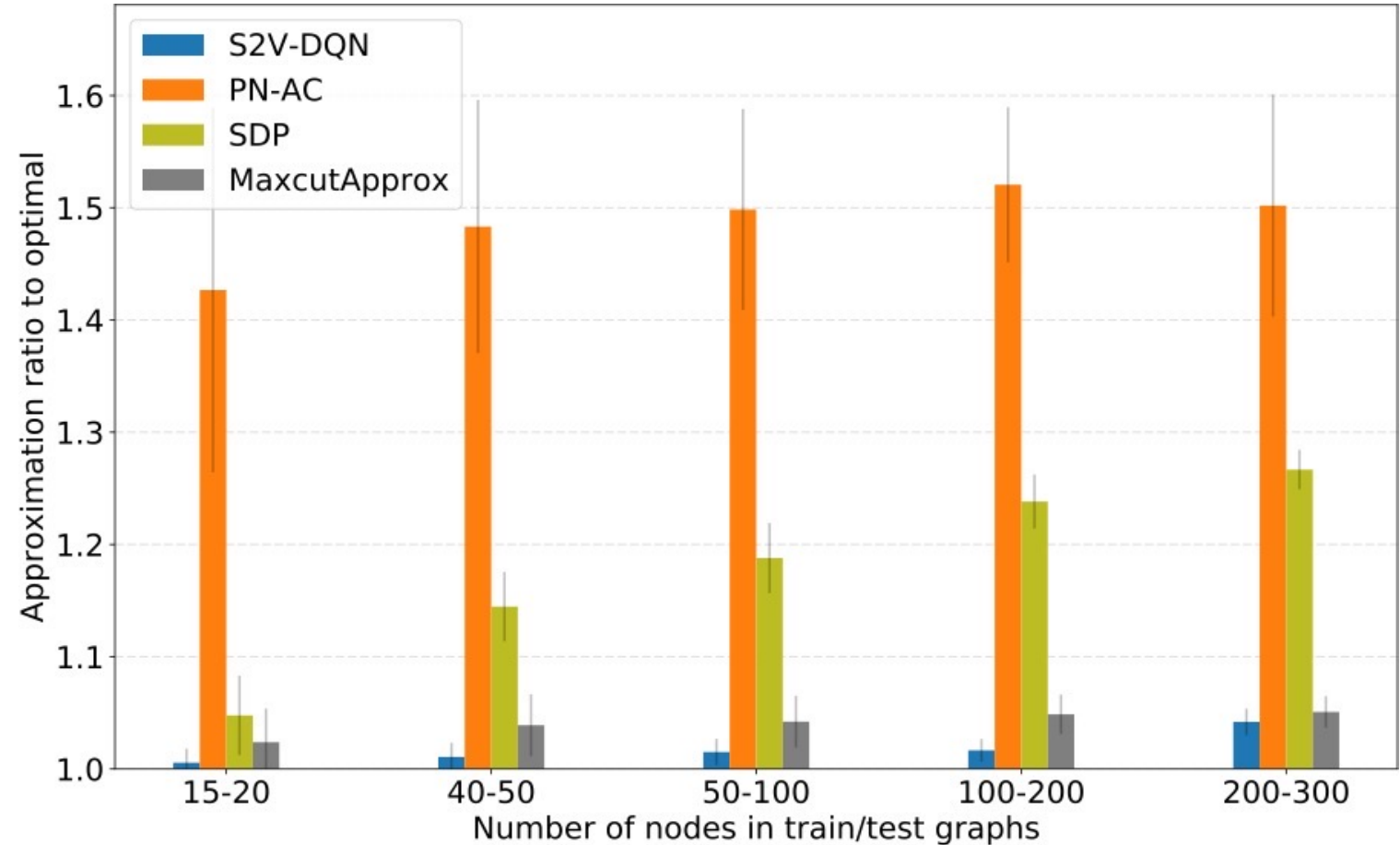
Barabasi-Albert
random graphs

Paper's approach

Another DL approach
[Bello et al., arXiv'16]

Goemans-Williamson
algorithm

Greedy algorithm
from first few slides



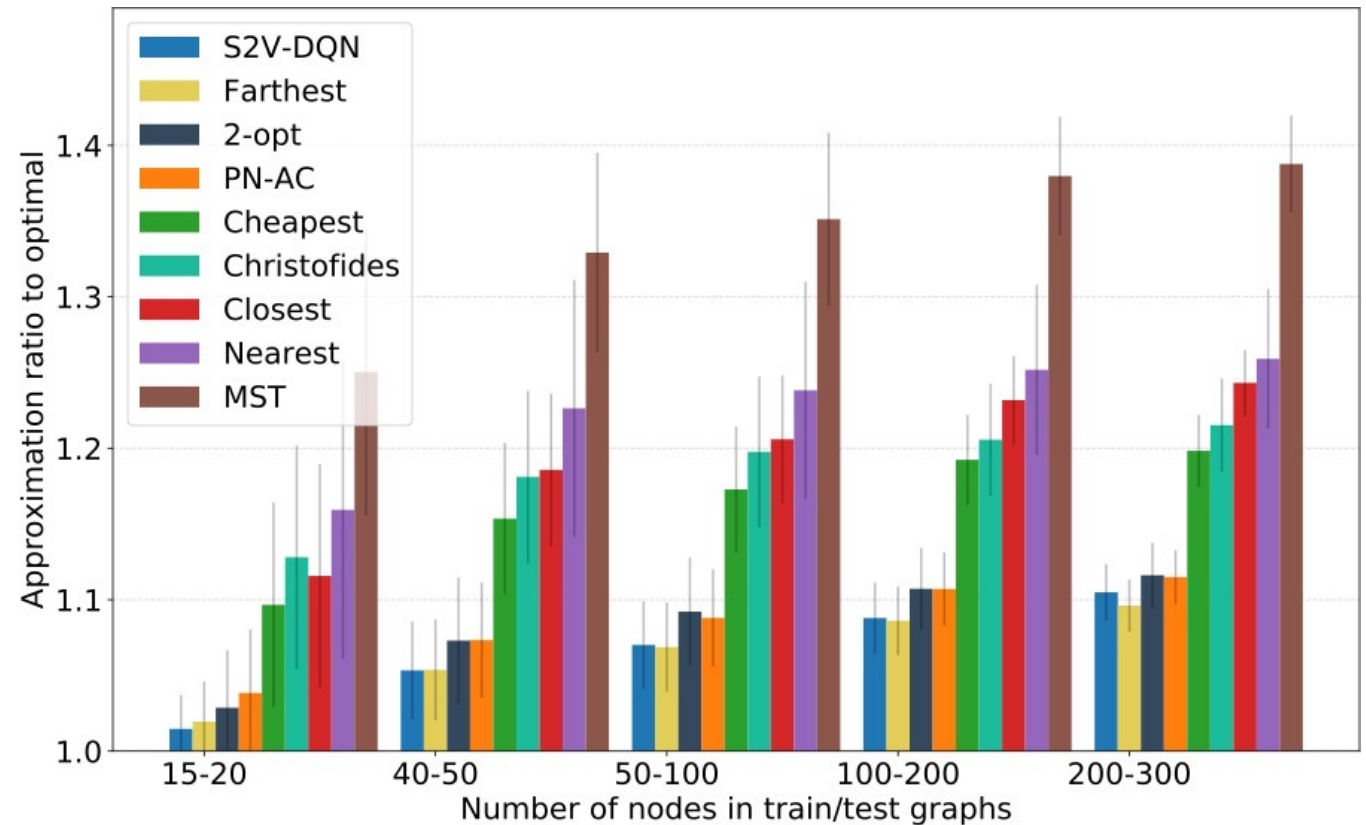
TSP

Uniform random points on 2-D grid

Paper's approach

- Initial subtour: 2 cities that are farthest apart
- Repeat the following:
 - Choose city that's *farthest* from any city in the subtour
 - Insert in position where it causes the smallest distance increase

[Rosenkrantz et al., SIAM JoC'77]



Runtime comparisons

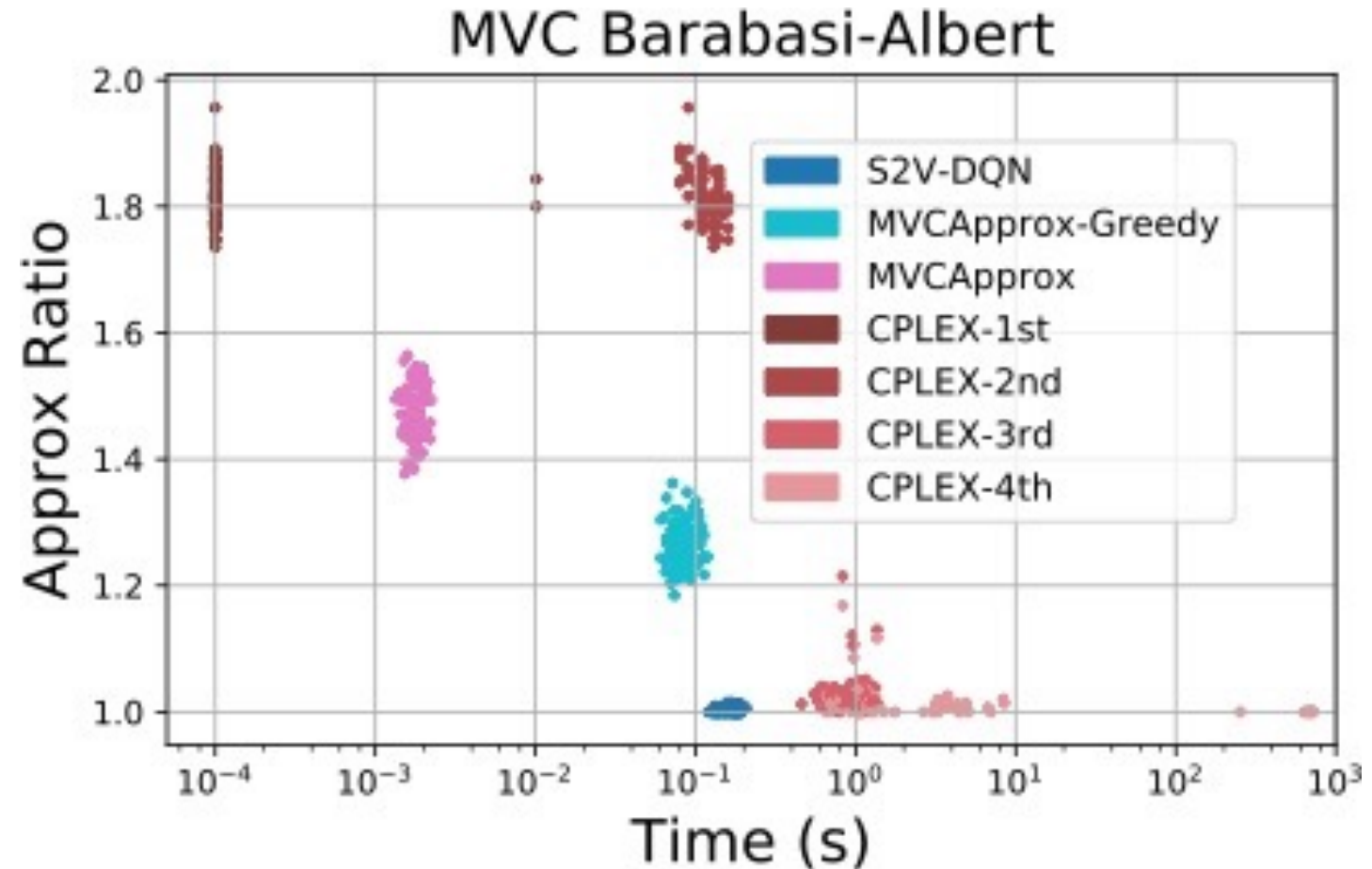
Paper's approach

Greedy algorithm from first few slides

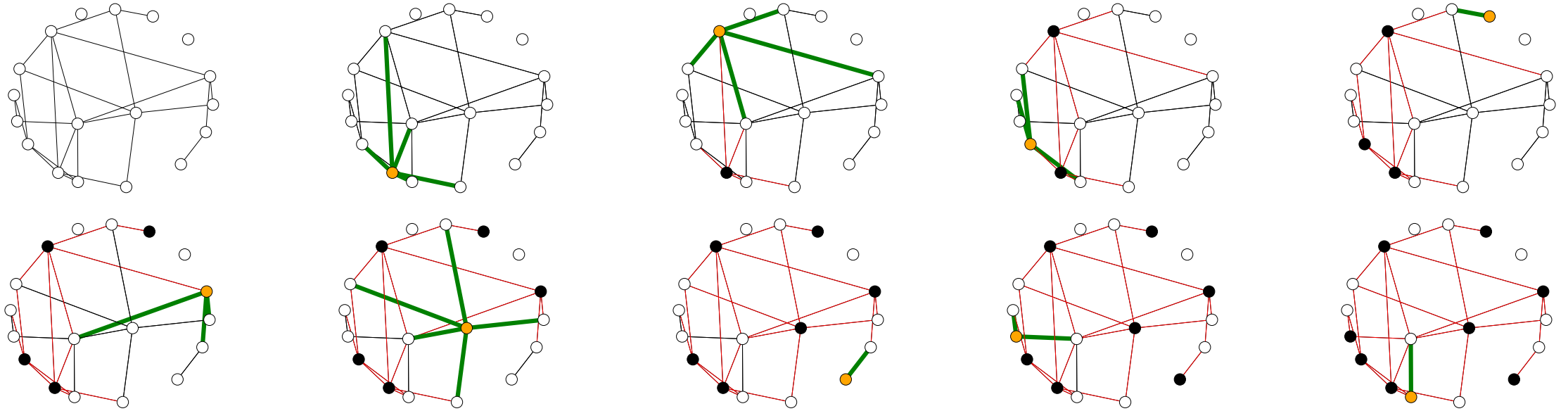
2-approximation algorithm

CPLEX-1st: 1st feasible solution found by CPLEX

CPLEX-2nd: 2nd feasible solution found by CPLEX



Min vertex cover visualization



Nodes seem to be selected to balance between:

- Degree
- Connectivity of the remaining graph

Summary

1 Applied techniques

- a. Graph neural networks
 - a. Neural algorithmic alignment
 - b. Variable selection for integer programming
- b. Learning greedy heuristics with RL

2 After the break: Theoretical guarantees

- a. Statistical guarantees for algorithm configuration
- b. Algorithms with predictions

Where much of my research has been

Summary

1 Applied techniques

- a. Graph neural networks
 - a. Neural algorithmic alignment
 - b. Variable selection for integer programming
- b. Learning greedy heuristics with RL

2 Theoretical guarantees

- a. **Statistical guarantees for algorithm configuration**
- b. Algorithms with predictions

Balcan, DeBlasio, Dick, Kingsford, Sandholm, **Vitercik**, STOC'21

Algorithm configuration

Example: **Integer programming solvers**

Most popular tool for solving combinatorial (& nonconvex) problems



Routing



Manufacturing



Scheduling



Planning



Finance

Algorithm configuration

IP solvers (CPLEX, Gurobi) have a **ton** parameters

- CPLEX has **170-page** manual describing **172** parameters
- Tuning by hand is notoriously **slow, tedious,** and **error-prone**

CPX_PARAM_NODEFILEIND 100	CPX_PARAM_TRELIM 160	CPX_PARAM_RANDOMSEED 130	CPXPARAM_MIP_Pool_RelGap 148	CPX_PARAM_FLOWCOVERS 70	CPX_PARAM_BRDIR 39
CPX_PARAM_NODELIM 101	CPX_PARAM_TUNINGDETILIM 160	CPX_PARAM_REDUCE 131	CPXPARAM_MIP_Pool_Replace 151	CPX_PARAM_FLOWPATHS 71	CPX_PARAM_BTTOL 40
CPX_PARAM_NODESEL 102	CPX_PARAM_TUNINGDISPLAY 162	CPX_PARAM_REINV 131	CPXPARAM_MIP_Strategy_Branch 39	CPX_PARAM_FPHEUR 72	CPX_PARAM_CALCQCPCDUALS 41
CPX_PARAM_NUMERICALEMPHASIS 102	CPX_PARAM_TUNINGMEASURE 163	CPX_PARAM_RELAXPREIND 132	CPXPARAM_MIP_Strategy_MIQCPStrat 93	CPX_PARAM_FRACCAND 73	CPX_PARAM_CLIQUES 42
CPX_PARAM_NZREADLIM 103	CPX_PARAM_TUNINGREPEAT 164	CPX_PARAM_RELOBJDIF 133	CPXPARAM_MIP_Strategy_StartAlgorithm 139	CPX_PARAM_FRACCUTS 73	CPX_PARAM_CLOCKTYPE 43
CPX_PARAM_OBJDIF 104	CPX_PARAM_TUNINGTILIM 165	CPX_PARAM_REPAIRTRIES 133	CPXPARAM_MIP_Strategy_VariableSelect 166	CPX_PARAM_FRACPASS 74	CPX_PARAM_CLONELOG 43
CPX_PARAM_OBJLLIM 105	CPX_PARAM_VARSEL 166	CPX_PARAM_REPEATPRESOLVE 134	CPXPARAM_MIP_SubMIP_NodeLimit 155	CPX_PARAM_GUBCOVERS 75	CPX_PARAM_COEREDIND 44
CPX_PARAM_OBJULIM 105	CPX_PARAM_WORKDIR 167	CPX_PARAM_RINSHEUR 135	CPXPARAM_OptimalityTarget 106	CPX_PARAM_HEURFREQ 76	CPX_PARAM_COLREADLIM 45
CPX_PARAM_PARALLELMODE 108	CPX_PARAM_WORKMEM 168	CPX_PARAM_RLT 136	CPXPARAM_Output_WriteLevel 169	CPX_PARAM_IMPLBD 76	CPX_PARAM_CONFLICTDISPLAY 46
CPX_PARAM_PERIND 110	CPX_PARAM_WRITELEVEL 169	CPX_PARAM_ROWREADLIM 141	CPXPARAM_Preprocessing_Aggregator 19	CPX_PARAM_INTSOLFILEPREFIX 78	CPX_PARAM_COVERS 47
CPX_PARAM_PERLIM 111	CPX_PARAM_ZEROHALFCUTS 170	CPX_PARAM_SCAIND 142	CPXPARAM_Preprocessing_Fill 19	CPX_PARAM_INTSOLLIM 79	CPX_PARAM_CPUMASK 48
CPX_PARAM_POLISHAFTERDETTIME 111	CPXPARAM_Benders_Strategy 30	CPX_PARAM_SCRIND 143	CPXPARAM_Preprocessing_Linear 120	CPX_PARAM_ITLIM 80	CPX_PARAM_CRAIN 50
CPX_PARAM_POLISHAFTEREPAGAP 112	CPXPARAM_Benders_Tolerances_feasibilitycut 35	CPX_PARAM_SIFTALG 143	CPXPARAM_Preprocessing_Reduce 131	CPX_PARAM_LANDPCUTS 82	CPX_PARAM_CUTLO 51
CPX_PARAM_POLISHAFTEREPGAP 113	CPXPARAM_Benders_Tolerances_optimalitycut 36	CPX_PARAM_SIFTDISPLAY 144	CPXPARAM_Preprocessing_Symmetry 156	CPX_PARAM_LBHEUR 81	CPX_PARAM_CUTPASS 52
CPX_PARAM_POLISHAFTERINTSOL 114	CPXPARAM_Conflict_Algorithm 46	CPX_PARAM_SIFTITLIM 145	CPXPARAM_Read_DataCheck 54	CPX_PARAM_LPMETHOD 136	CPX_PARAM_CUTSFACTOR 52
CPX_PARAM_POLISHAFTERNODE 115	CPXPARAM_CPUmask 48	CPX_PARAM_SIMDISPLAY 145	CPXPARAM_Read_Scale 142	CPX_PARAM_MCFCUTS 82	CPX_PARAM_CUTUP 53
CPX_PARAM_POLISHAFTERTIME 116	CPXPARAM_DistMIP_Rampup_Duration 128	CPX_PARAM_SINGLIM 146	CPXPARAM_ScreenOutput 143	CPX_PARAM_MEMORYEMPHASIS 83	CPXPARAM_DATACHECK 54
CPX_PARAM_POLISHTIME (deprecated) 116	CPXPARAM_LPMethod 136	CPX_PARAM_SOLNPOOLGAP 146	CPXPARAM_Sifting_Algorithm 143	CPX_PARAM_MIPCBREDLP 84	CPX_PARAM_DEPIND 55
CPX_PARAM_POPULATELIM 117	CPXPARAM_MIP_Cuts_BQP 38	CPX_PARAM_SOLNPOOLCAPACITY 147	CPXPARAM_Sifting_Display 144	CPX_PARAM_MIPDISPLAY 85	CPX_PARAM_DETTILIM 56
CPX_PARAM_PPRIND 118	CPXPARAM_MIP_Cuts_LocallyImplied 77	CPX_PARAM_SOLNPOOLREPLACE 151	CPXPARAM_Sifting_Iterations 145	CPX_PARAM_MIPEMPHASIS 87	CPX_PARAM_DISJCUTS 57
CPX_PARAM_PREDUAL 119	CPXPARAM_MIP_Cuts_RLT 136	CPX_PARAM_SOLUTIONTARGET (deprecated: see CPXPARAM_OptimalityTarget 106)	CPXPARAM_Simplex_Display 145	CPX_PARAM_MIPINTERVAL 88	CPX_PARAM_DIVETYPE 58
CPX_PARAM_PREIND 120	CPXPARAM_MIP_Cuts_ZeroHalfCut 170	CPX_PARAM_THREADS 157	CPXPARAM_Simplex_Limits_Singularity 146	CPX_PARAM_MIPKAPPASTATS 89	CPX_PARAM_DPRIIND 59
CPX_PARAM_PRLINEAR 120	CPXPARAM_MIP_Limits_CutsFactor 52	CPXPARAM_TimeLimit 159	CPXPARAM_SolutionType 152	CPX_PARAM_MIPORDIND 90	CPX_PARAM_EACHCUTLIM 60
CPX_PARAM_PREPASS 121	CPXPARAM_MIP_Limits_RampupDetTimeLimit 127	CPXPARAM_Tune_DefTimeLimit 160	CPXPARAM_Threads 157	CPX_PARAM_MIPORDTYPE 91	CPX_PARAM_EPAGAP 61
CPX_PARAM_PRESLVND 122	CPXPARAM_MIP_Limits_RampupTimeLimit 128	CPXPARAM_Tune_Display 162	CPXPARAM_TimeLimit 159	CPX_PARAM_MIPSEARCH 92	CPX_PARAM_EPGAP 61
CPX_PARAM_PRICELIM 123	CPXPARAM_MIP_Limits_Solutions 79	CPXPARAM_Tune_Measure 163	CPXPARAM_Tune_DefTimeLimit 160	CPX_PARAM_MIQCPSTRAT 93	CPX_PARAM_EPINT 62
CPX_PARAM_PROBE 123	CPXPARAM_MIP_Limits_StrongCand 154	CPXPARAM_Tune_Repeat 164	CPXPARAM_Tune_Display 162	CPX_PARAM_MIRCUTS 94	CPX_PARAM_EPMRK 64
CPX_PARAM_PROBEDETTIME 124	CPXPARAM_MIP_Limits_StrongIt 154	CPXPARAM_Tune_Measure 163	CPXPARAM_Tune_Measure 163	CPX_PARAM_MPSLONGNUM 94	CPX_PARAM_EPOPT 65
CPX_PARAM_PROBETIME 124	CPXPARAM_MIP_Limits_TreeMemory 160	CPXPARAM_Tune_Repeat 164	CPXPARAM_Tune_Repeat 164	CPX_PARAM_NETDISPLAY 95	CPX_PARAM_EPPER 65
CPX_PARAM_QPMAKEPSDIND 125	CPXPARAM_MIP_OrderType 91	CPXPARAM_Tune_SubALG 165	CPXPARAM_Tune_SubALG 165	CPX_PARAM_NETEPOPT 96	CPX_PARAM_EPRELAX 66
CPX_PARAM_QPMETHOD 138	CPXPARAM_MIP_Pool_AbsGap 146	CPXPARAM_WorkDir 167	CPXPARAM_WorkDir 167	CPX_PARAM_NETEPRHS 96	CPX_PARAM_EPRHS 67
CPX_PARAM_QPNZREADLIM 126	CPXPARAM_MIP_Pool_Capacity 147	CPXPARAM_WorkMem 168	CPXPARAM_WorkMem 168	CPX_PARAM_NETFIND 97	CPX_PARAM_FEASOPTMODE 68
	CPXPARAM_MIP_Pool_Intensity 149	CraInd 50		CPX_PARAM_NETITLIM 98	CPX_PARAM_FILEENCODING 69
		CPX_PARAM_TILIM 159		CPX_PARAM_NETPPRIIND 98	

Algorithm configuration

IP solvers (CPLEX, Gurobi) have a **ton** parameters

- CPLEX has **170-page** manual describing **172** parameters
- Tuning by hand is notoriously **slow, tedious**, and **error-prone**

What's the best **configuration** for the application at hand?



Best configuration for **routing** problems
likely not suited for **scheduling**



Running example: Sequence alignment

Goal: Line up pairs of strings

Applications: Biology, natural language processing, etc.



Did you mean: [vitercik](#)

Sequence alignment algorithms

Input: Two sequences S and S'

Output: Alignment of S and S'

$S = A C T G$
 $S' = G T C A$

Gap
↓
A - - C T G
- G T C A -
↑ ↑ ↑
Insertion/deletion (*indel*) Match Mismatch

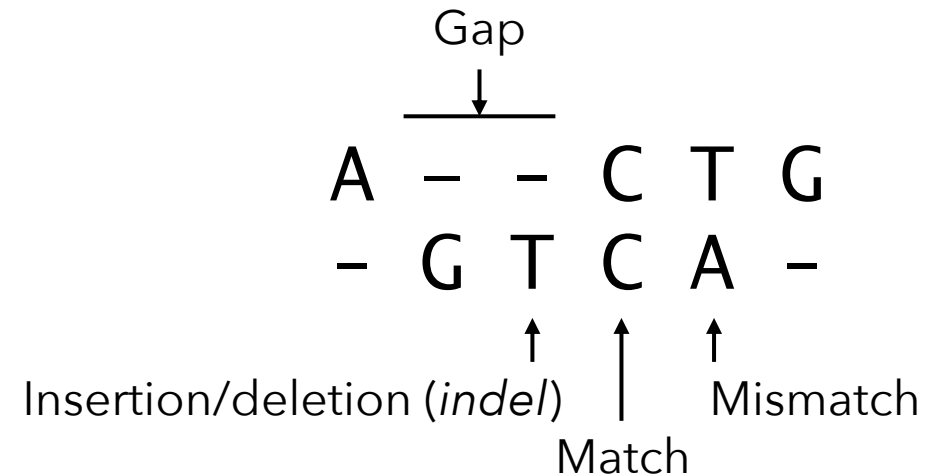
Sequence alignment algorithms

Standard algorithm with parameters $\rho_1, \rho_2, \rho_3 \geq 0$:

Return alignment maximizing:

$$(\# \text{ matches}) - \rho_1 \cdot (\# \text{ mismatches}) - \rho_2 \cdot (\# \text{ indels}) - \rho_3 \cdot (\# \text{ gaps})$$

$S = A C T G$
 $S' = G T C A$



Sequence alignment algorithms


Can sometimes access **ground-truth, reference** alignment

E.g., in computational biology: Bahr et al., Nucleic Acids Res.'01; Raghava et al., BMC Bioinformatics '03; Edgar, Nucleic Acids Res.'04; Walle et al., Bioinformatics'04

Requires extensive manual alignments
...rather just run parameterized algorithm

How to tune algorithm's parameters?

*"There is **considerable disagreement** among molecular biologists about the **correct choice**" [Gusfield et al. '94]*



A	-	-	C	T	G
-	G	T	C	A	-

Sequence alignment algorithms

-GRTCPKPDDLPFSTVVP-LKTFYEPGEEITYSCKPGYVSRGGMRKFICPLTGLWPINTLKCTP
E-VKCPFPSRPDNGFVNYPKPTLYYKDKATFGCHDGYSLDGP-EEIECTKLGNEWSAMPSC-KA

Ground-truth alignment of protein sequences

Sequence alignment algorithms

-GRTCPKPDDLPFSTVVP-LKTFYEPEEEITYSCKPGYVSRGGMRKFICPLTGLWPINTLKCTP
E-VKCPFPSRPDNGFVNYPKPTLYYKDKATFGCHDGYSLDGP-EEIECTKLGNSAMPSC-KA

Ground-truth alignment of protein sequences

GRTCP---KPDDLPFSTVVPLKTFYEPEEEITYSCKPGYVSRGGMRKFICPLTGLWPINTLKCTP
EVKCPFPSRPDN-GFVNYPKPTLYYK-DKATFGCHDGY-SLDGPEEIECTKLGNS-AMPSCKA

Alignment by algorithm with **poorly-tuned** parameters

Sequence alignment algorithms

-GRTCPKPDDL PFSTVVP-LKTFYEPGEEITYSCKPGYVSRGGMRKFICPLTGLWPINTLKCTP
E-VKCPFPSRPDNGFVNYP AKPTLYYKDKATFGCHDGYSLDGP-EEIECTKLGNSAMPSC-KA

Ground-truth alignment of protein sequences

GRTCP---KPDDL PFSTVVPLKTFYEPGEEITYSCKPGYVSRGGMRKFICPLTGLWPINTLKCTP
EVKCPFPSRPDN-GFVNYP AKPTLYYK-DKATFGCHDGY-SLDGPEEIECTKLGNS-AMPSCKA

Alignment by algorithm with **poorly-tuned** parameters

GRTCPKPDDL PFSTV-VPLKTFYEPGEEITYSCKPGYVSRGGMRKFICPLTGLWPINTLKCTP
EVKCPFPSRPDNGFVNYP AKPTLYYKDKATFGCHDGY-SLDGPEEIECTKLGNSA-MPSCKA

Alignment by algorithm with **well-tuned** parameters

Automated parameter tuning procedure

1. Fix parameterized algorithm
2. Receive training set T of "typical" inputs



3. Find parameter setting w/ good avg performance over T
Runtime, solution quality, etc.

Automated parameter tuning procedure

1. Fix parameterized algorithm
2. Receive training set T of "typical" inputs



3. Find parameter setting w/ good avg performance over T

On average, output alignment is close to reference alignment

Automated parameter tuning procedure

1. Fix parameterized algorithm
2. Receive training set T of "typical" inputs



3. Find parameter setting w/ good avg performance over T

Key question:

How to find parameter setting with good avg performance?

Automated parameter tuning procedure

Key question:

How to find parameter setting with good avg performance?



E.g., for sequence alignment:
algorithm by Gusfield et al. ['94]

Many other generic search strategies

E.g., Hutter et al. [JAIR'09, LION'11], Ansótegui et al. [CP'09], ...

Automated parameter tuning procedure

1. Fix parameterized algorithm
2. Receive training set T of "typical" inputs

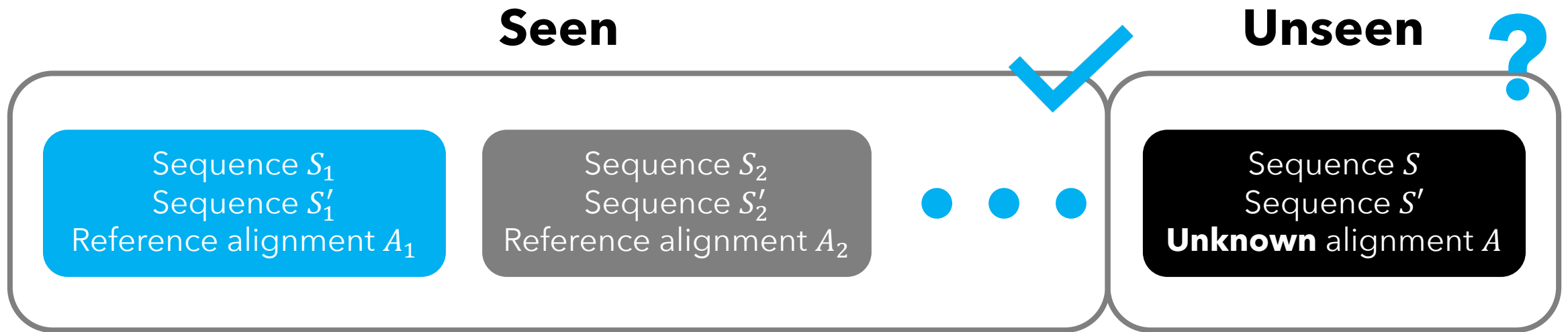


3. Find parameter setting w/ good avg performance over T

Key question (focus of this section):

Will that parameter setting have good **future** performance?

Automated parameter tuning procedure



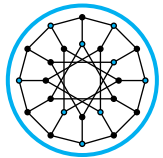
Key question (focus of this section):

Will that parameter setting have good **future** performance?

Generalization

Key question (focus of this section):

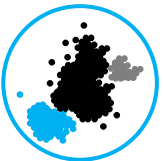
Good performance on **average** over **training set** implies good **future** performance?



Greedy algorithms

Gupta, Roughgarden, ITCS'16

First to ask question for algorithm configuration



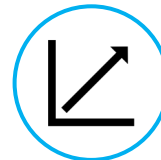
Clustering

Balcan, Nagarajan, **V**, White, COLT'17
Garg, Kalai, NeurIPS'18
Balcan, Dick, White, NeurIPS'18
Balcan, Dick, Lang, ICLR'20



Search

Sakaue, Oki, NeurIPS'22



Numerical linear algebra

Bartlett et al., COLT'22

And many other areas...

This section: Main result

Key question (focus of this section):

Good performance on **average** over **training set** implies good **future** performance?

Answer this question for any parameterized algorithm where:

Performance is **piecewise-structured** function of parameters

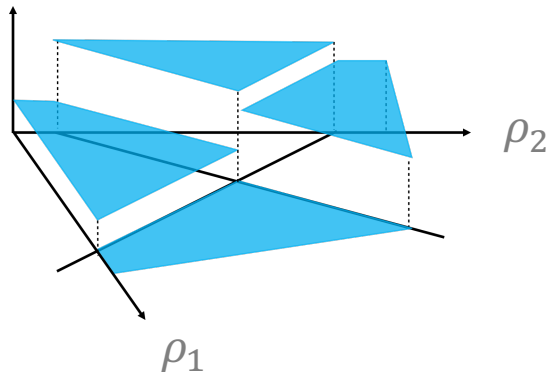
Piecewise constant, linear, quadratic, ...

This section: Main result

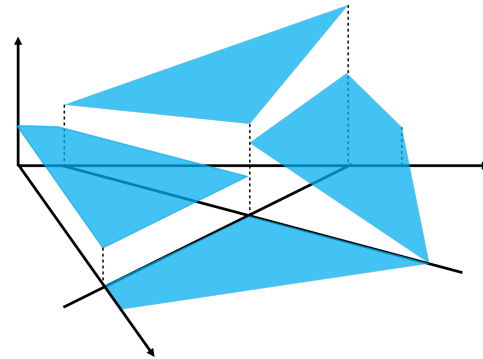
Performance is **piecewise-structured** function of parameters

Piecewise constant, linear, quadratic, ...

Algorithmic
performance
on fixed input



Piecewise constant



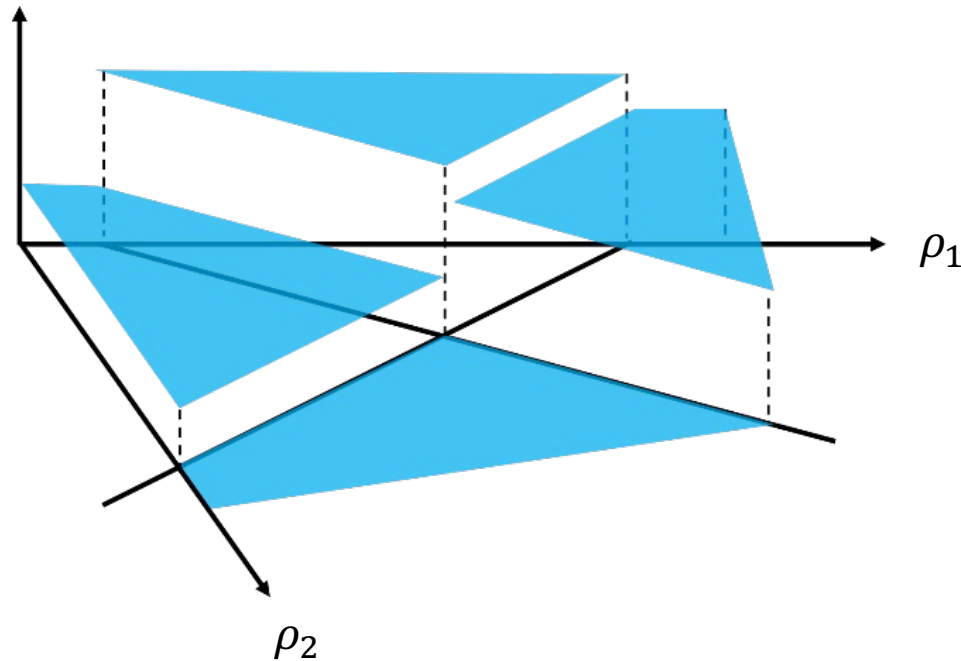
Piecewise linear



Piecewise ...

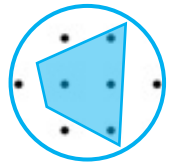
Example: Sequence alignment

Distance between **algorithm's output** given S, S'
and **ground-truth** alignment is p-wise constant



Piecewise structure

Piecewise structure unifies **seemingly disparate** problems:



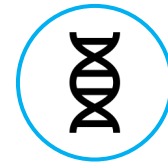
Integer programming

Balcan, Prasad, Sandholm, **V**, NeurIPS'21
Balcan, Prasad, Sandholm, **V**, NeurIPS'22
Balcan, Dick, Sandholm, **V**, JACM'24



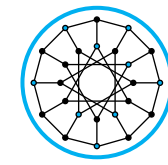
Clustering

Balcan, Nagarajan, **V**, White, COLT'17
Balcan, Dick, White, NeurIPS'18
Balcan, Dick, Lang, ICLR'20



Computational biology

Balcan, DeBlasio, Dick, Kingsford,
Sandholm, **V**, STOC'21



Greedy algorithms

Gupta, Roughgarden, ITCS'16



Mechanism configuration

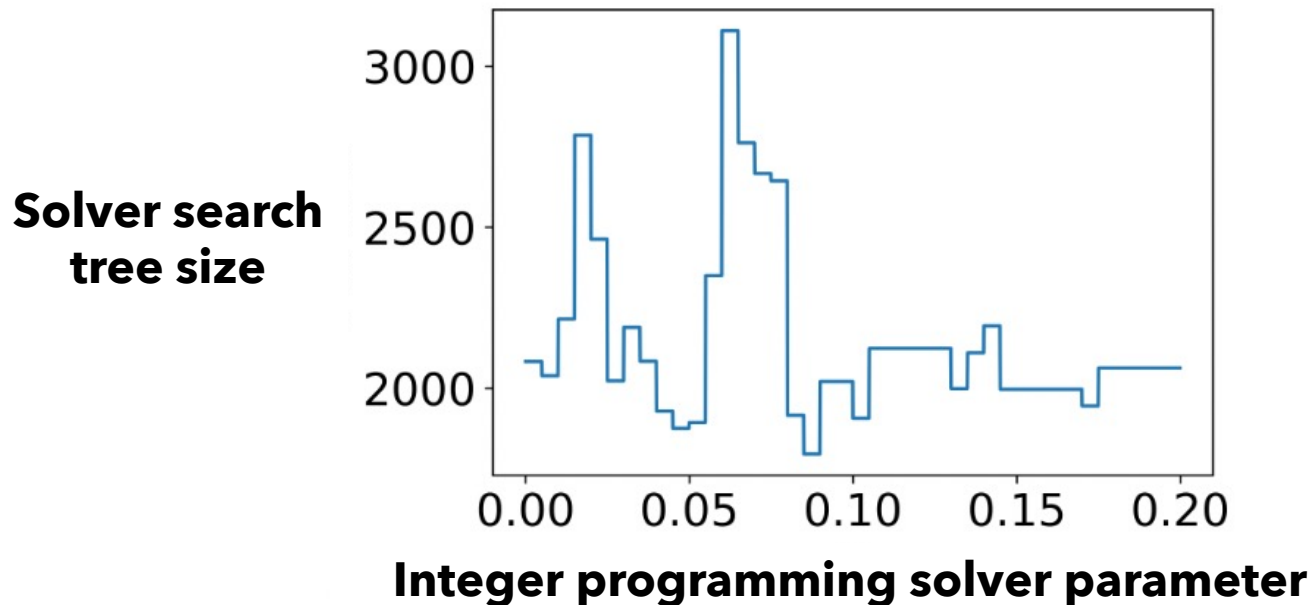
Balcan, Sandholm, **V**, OR'24

Ties to a long line of research on machine learning for **revenue maximization**

Likhodedov, Sandholm, AAAI'04, '05; Balcan, Blum, Hartline, Mansour, FOCS'05; Elkind, SODA'07;
Cole, Roughgarden, STOC'14; Mohri, Medina, ICML'14; Devanur, Huang, Psomas, STOC'16; ...

Primary challenge

Algorithmic performance is a **volatile** function of parameters
Complex connection between parameters and performance



Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
 - i. **Model**
 - ii. Piecewise-structured algorithmic performance
 - iii. Main result
 - iv. Application: Sequence alignment
 - v. Online algorithm configuration
2. Algorithms with predictions

Model

\mathbb{R}^d : Set of all parameter settings

\mathcal{X} : Set of all inputs

Example: Sequence alignment

\mathbb{R}^3 : Set of alignment algorithm parameter settings

\mathcal{X} : Set of sequence pairs



$S = A C T G$
 $S' = G T C A$

One sequence pair $x = (S, S') \in \mathcal{X}$

Algorithmic performance

$u_{\rho}(x)$ = utility of algorithm parameterized by $\rho \in \mathbb{R}^d$ on input x
E.g., runtime, solution quality, distance to ground truth, ...

Assume $u_{\rho}(x) \in [-1, 1]$

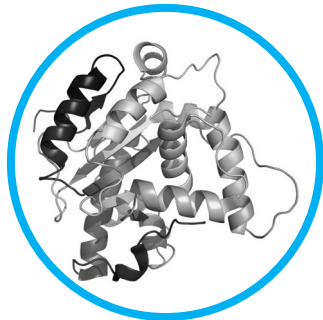
Can be generalized to $u_{\rho}(x) \in [-H, H]$

Model

Standard assumption: Unknown distribution \mathcal{D} over inputs
Distribution models specific application domain at hand



E.g., distribution over pairs of DNA strands



E.g., distribution over pairs of protein sequences

Generalization bounds

Key question: For any parameter setting ρ ,
is **average** utility on training set close to **expected** utility?
future



Generalization bounds

Key question: For any parameter setting ρ ,
is **average** utility on training set close to **expected** utility?

Formally: Given samples $x_1, \dots, x_N \sim \mathcal{D}$, for any ρ ,

$$\left| \underbrace{\frac{1}{N} \sum_{i=1}^N u_{\rho}(x_i)}_{\text{Empirical average utility}} - \underbrace{\mathbb{E}_{x \sim \mathcal{D}}[u_{\rho}(x)]}_{\text{Expected utility}} \right| \leq ?$$

Good **average empirical** utility \rightarrow Good **expected** utility

Outline (theoretical guarantees)

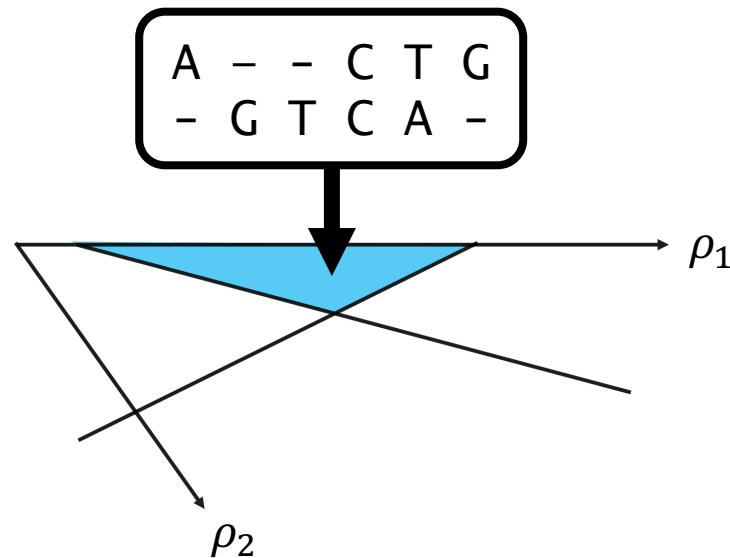
1. Statistical guarantees for algorithm configuration
 - i. Model
 - ii. Piecewise-structured algorithmic performance**
 - a. Example: Sequence alignment**
 - b. Dual function definition
 - iii. Main result
 - iv. Application: Sequence alignment
 - v. Online algorithm configuration
2. Algorithms with predictions

Sequence alignment algorithms

Lemma:

For any pair S, S' , there's a partition of \mathbb{R}^3 s.t. in any region, algorithm's output is fixed across all parameters in region

$S = A C T G$
 $S' = G T C A$

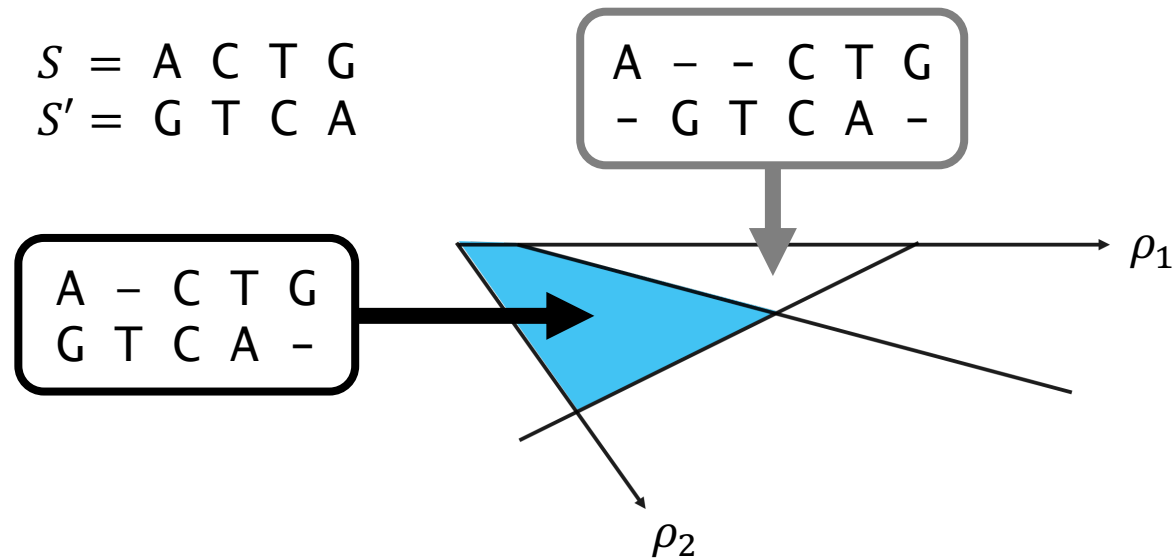


Sequence alignment algorithms

Lemma:

Defined by $(\max\{|S|, |S'|\})^3$ hyperplanes

For any pair S, S' , there's a partition of \mathbb{R}^3 s.t. in any region, algorithm's output is fixed across all parameters in region

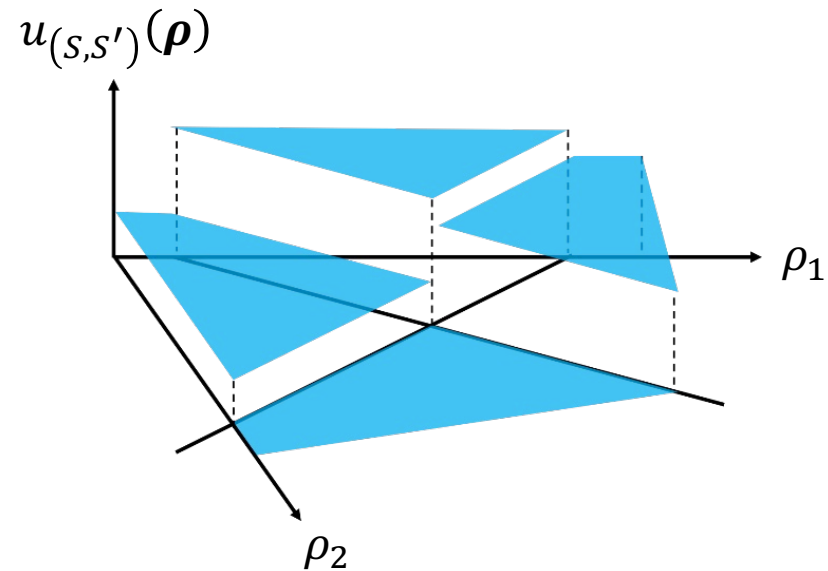


Piecewise-constant utility function

Corollary:

Utility is piecewise constant function of parameters

Distance between algorithm's output and ground-truth alignment



Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
 - i. Model
 - ii. Piecewise-structured algorithmic performance
 - a. Example: Sequence alignment
 - b. Dual function definition**
 - iii. Main result
 - iv. Application: Sequence alignment
 - v. Online algorithm configuration
2. Algorithms with predictions

Primal & dual classes

$u_{\rho}(x)$ = utility of algorithm parameterized by $\rho \in \mathbb{R}^d$ on input x
 $\mathcal{U} = \{u_{\rho}: \mathcal{X} \rightarrow \mathbb{R} \mid \rho \in \mathbb{R}^d\}$ “Primal” function class

Typically, prove guarantees by bounding **complexity** of \mathcal{U}

Challenge: \mathcal{U} is gnarly

E.g., in sequence alignment:

- Each domain element is a pair of sequences
- Unclear how to plot or visualize functions u_{ρ}
- No obvious notions of Lipschitz continuity or smoothness to rely on

Primal & dual classes

$u_{\rho}(x)$ = utility of algorithm parameterized by $\rho \in \mathbb{R}^d$ on input x
 $\mathcal{U} = \{u_{\rho}: \mathcal{X} \rightarrow \mathbb{R} \mid \rho \in \mathbb{R}^d\}$ "Primal" function class

$u_x^*(\rho)$ = utility as function of parameters

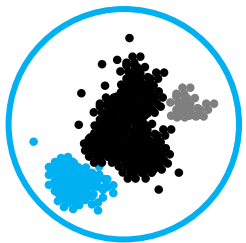
$$u_x^*(\rho) = u_{\rho}(x)$$

$\mathcal{U}^* = \{u_x^*: \mathbb{R}^d \rightarrow \mathbb{R} \mid x \in \mathcal{X}\}$ "Dual" function class

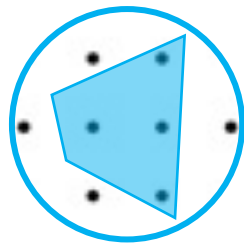
- Dual functions have simple, Euclidean domain
- Often have ample structure can use to bound complexity of \mathcal{U}

Piecewise-structured functions

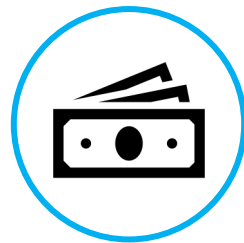
Dual functions $u_x^*: \mathbb{R}^d \rightarrow \mathbb{R}$ are **piecewise-structured**



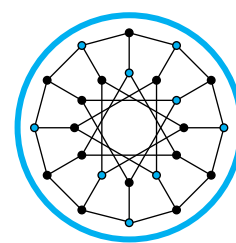
Clustering
algorithm
configuration



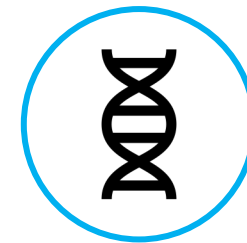
Integer programming
algorithm
configuration



Selling mechanism
configuration



Greedy
algorithm
configuration



Computational biology
algorithm
configuration



Voting mechanism
configuration

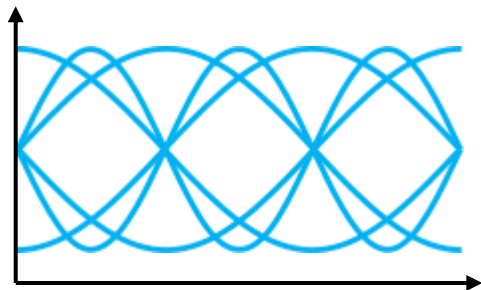
Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
 - i. Model
 - ii. Piecewise-structured algorithmic performance
 - iii. Main result**
 - iv. Application: Sequence alignment
 - v. Online algorithm configuration
2. Algorithms with predictions

Intrinsic complexity

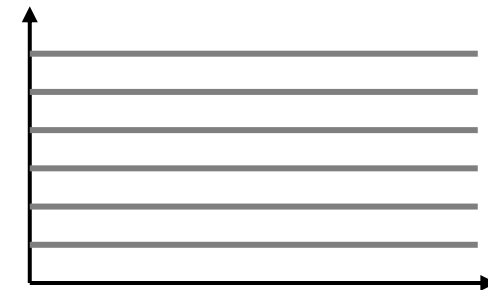
“Intrinsic complexity” of function class \mathcal{G}

- Measures how well functions in \mathcal{G} fit complex patterns
- Specific ways to quantify “intrinsic complexity”:
 - VC dimension
 - Pseudo-dimension



More complex

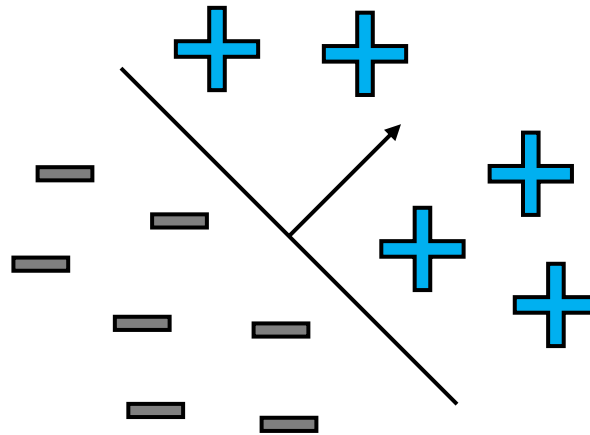
Less complex



VC dimension

Complexity measure for binary-valued function classes \mathcal{F}
(Classes of functions $f: \mathcal{Y} \rightarrow \{-1, 1\}$)

E.g., linear separators



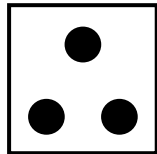
VC dimension

Size of the largest set $\mathcal{S} \subseteq \mathcal{Y}$

that can be labeled in all $2^{|\mathcal{S}|}$ ways by functions in \mathcal{F}

Example: \mathcal{F} = Linear separators in \mathbb{R}^2

$$\text{VCdim}(\mathcal{F}) \geq 3$$



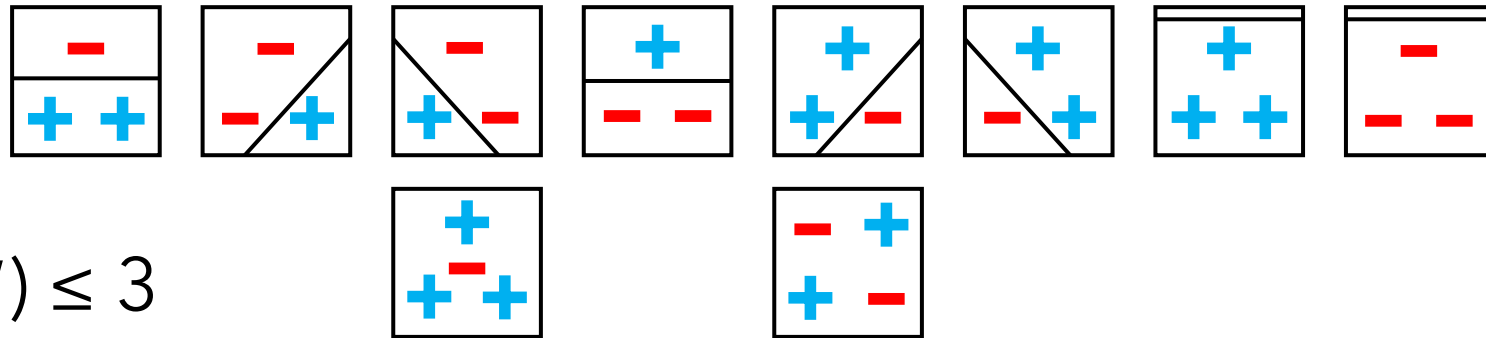
VC dimension

Size of the largest set $\mathcal{S} \subseteq \mathcal{Y}$

that can be labeled in all $2^{|\mathcal{S}|}$ ways by functions in \mathcal{F}

Example: \mathcal{F} = Linear separators in \mathbb{R}^2

$\text{VCdim}(\mathcal{F}) \geq 3$



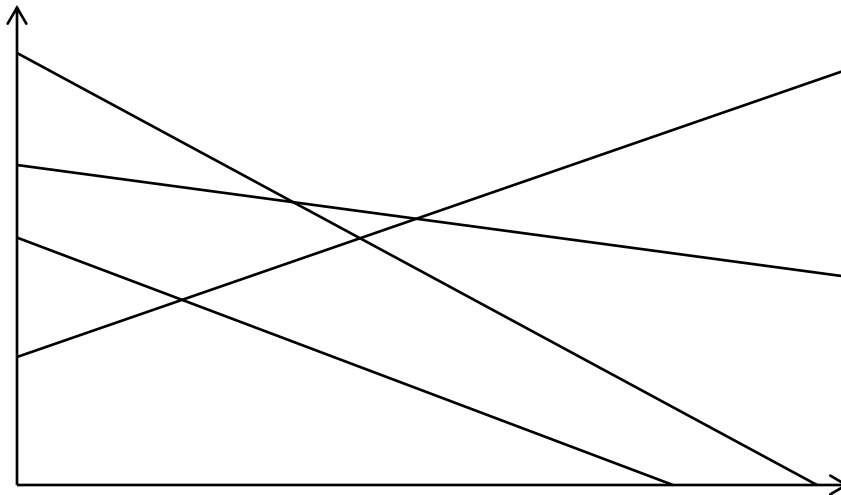
$\text{VCdim}(\mathcal{F}) \leq 3$

$\text{VCdim}(\{\text{Linear separators in } \mathbb{R}^d\}) = d + 1$

Pseudo-dimension

Complexity measure for real-valued function classes \mathcal{G}
(Classes of functions $g: \mathcal{Y} \rightarrow [-1,1]$)

E.g., affine functions



Pseudo-dimension of \mathcal{G}

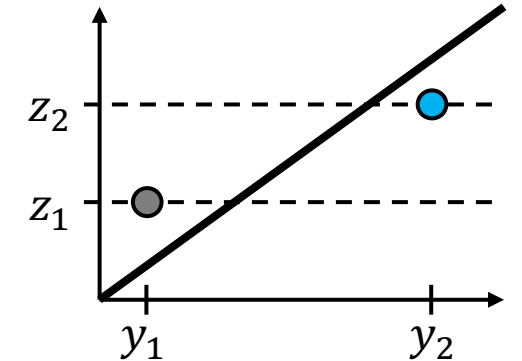
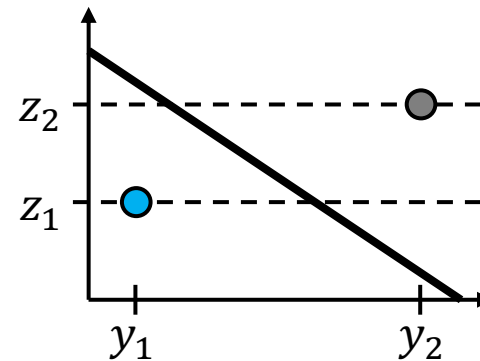
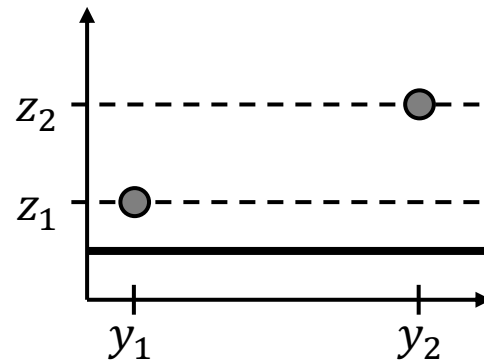
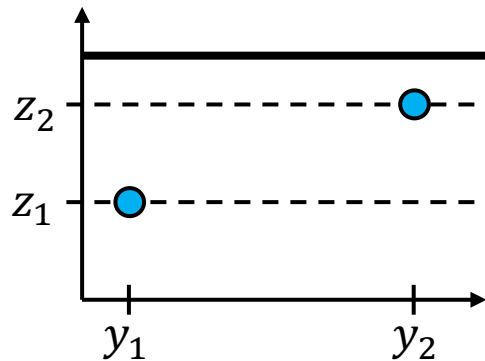
Size of the largest set $\{y_1, \dots, y_N\} \subseteq \mathcal{Y}$ s.t.:

for some *targets* $z_1, \dots, z_N \in \mathbb{R}$,

all 2^N above/below patterns achieved by functions in \mathcal{G}

Example: \mathcal{G} = Affine functions in \mathbb{R}

$\text{Pdim}(\mathcal{G}) \geq 2$



Can also show that $\text{Pdim}(\mathcal{G}) \leq 2$

Sample complexity using pseudo-dim

In the context of **algorithm configuration**:

- $\mathcal{U} = \{u_\rho : \rho \in \mathbb{R}^d\}$ measure algorithm **performance**
- For $\epsilon, \delta \in (0,1)$, let $N = O\left(\frac{\text{Pdim}(\mathcal{U})}{\epsilon^2} \log \frac{1}{\delta}\right)$
- With probability at least $1 - \delta$ over $x_1, \dots, x_N \sim \mathcal{D}, \forall \rho \in \mathbb{R}^d$,

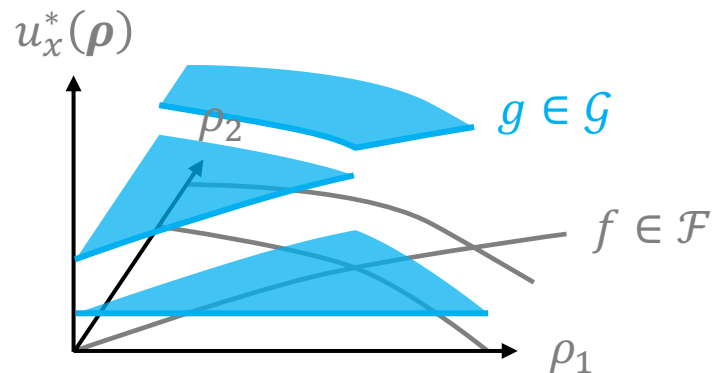
$$\left| \underbrace{\frac{1}{N} \sum_{i=1}^N u_\rho(x_i)}_{\text{Empirical average utility}} - \underbrace{\mathbb{E}_{x \sim \mathcal{D}}[u_\rho(x)]}_{\text{Expected utility}} \right| \leq \epsilon$$



Main result (informal)

Boundary functions $f_1, \dots, f_k \in \mathcal{F}$ partition \mathbb{R}^d s.t. in each region, $u_x^*(\boldsymbol{\rho}) = g(\boldsymbol{\rho})$ for some $g \in \mathcal{G}$.

Training set of size $\tilde{O}\left(\frac{\text{Pdim}(\mathcal{G}^*) + \text{VCdim}(\mathcal{F}^*) \log k}{\epsilon^2}\right)$ implies
WHP $\forall \boldsymbol{\rho}, |\text{avg utility over training set} - \text{exp utility}| \leq \epsilon$



Main result (informal)

Boundary functions $f_1, \dots, f_k \in \mathcal{F}$ partition \mathbb{R}^d s.t. in each region, $u_x^*(\boldsymbol{\rho}) = g(\boldsymbol{\rho})$ for some $g \in \mathcal{G}$.

Training set of size $\tilde{O}\left(\frac{\text{Pdim}(\mathcal{G}^*) + \text{VCdim}(\mathcal{F}^*) \log k}{\epsilon^2}\right)$ implies
WHP $\forall \boldsymbol{\rho}, |\text{avg utility over training set} - \text{exp utility}| \leq \epsilon$

\mathcal{F}, \mathcal{G} are typically very well structured

- \mathcal{G} = set of all **constant** functions $\Rightarrow \text{Pdim}(\mathcal{G}^*) = O(1)$
- \mathcal{G} = set of all **linear** functions in \mathbb{R}^d $\Rightarrow \text{Pdim}(\mathcal{G}^*) = O(d)$

Main result (informal)

Boundary functions $f_1, \dots, f_k \in \mathcal{F}$ partition \mathbb{R}^d s.t. in each region, $u_x^*(\boldsymbol{\rho}) = g(\boldsymbol{\rho})$ for some $g \in \mathcal{G}$.

Theorem:

$$\text{Pdim}(\mathcal{U}) = \tilde{O}(\text{Pdim}(\mathcal{G}^*) + \text{VCdim}(\mathcal{F}^*) \log k)$$

↑
Primal function class $\mathcal{U} = \{u_\rho \mid \rho \in \mathbb{R}^d\}$

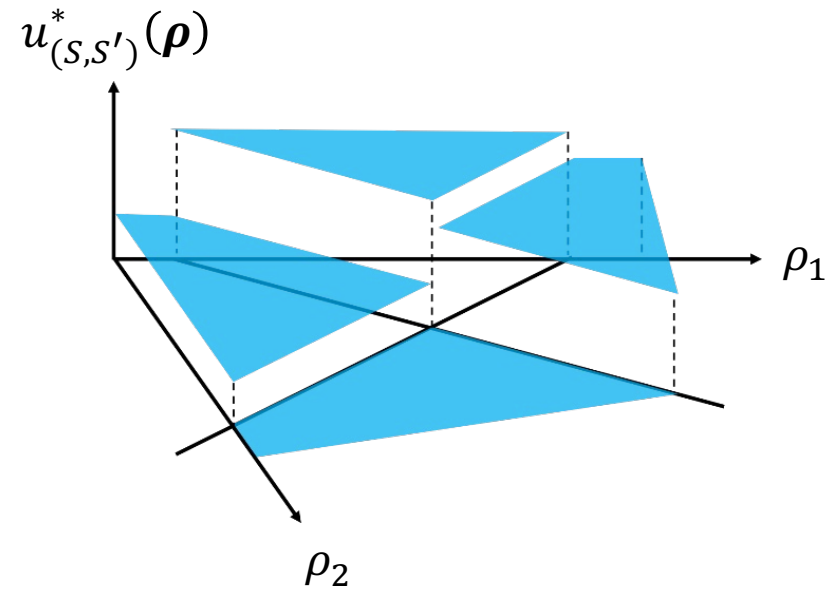
Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
 - i. Model
 - ii. Piecewise-structured algorithmic performance
 - iii. Main result
 - iv. Application: Sequence alignment**
 - v. Online algorithm configuration
2. Algorithms with predictions

Piecewise constant dual functions

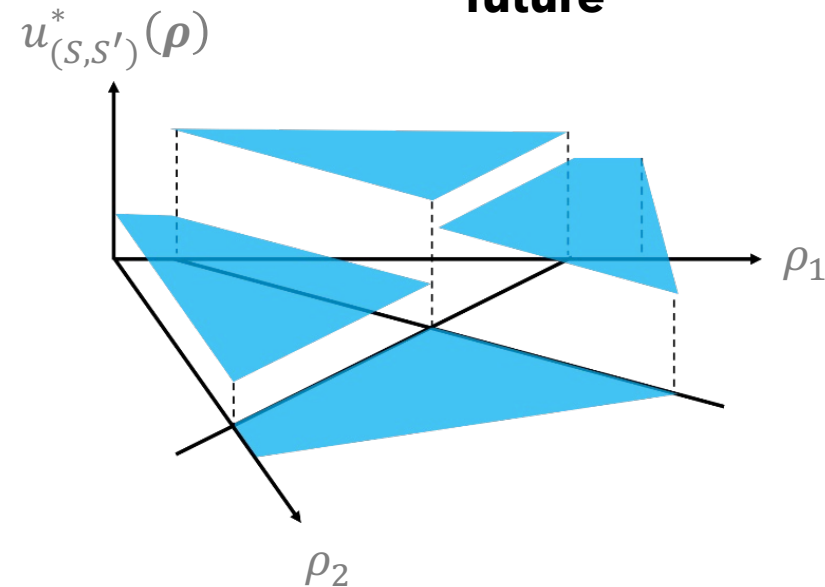
Lemma:

Utility is piecewise constant function of parameters

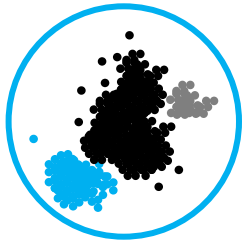


Sequence alignment guarantees

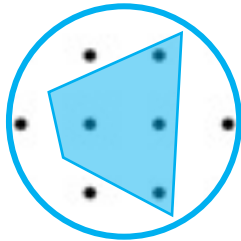
Theorem: Training set of size $\tilde{O}\left(\frac{\log(\text{seq. length})}{\epsilon^2}\right)$ implies WHP $\forall \rho$,
|**avg** utility over training set - **exp** utility| $\leq \epsilon$
future



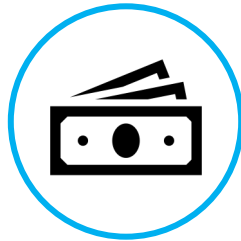
Many more applications



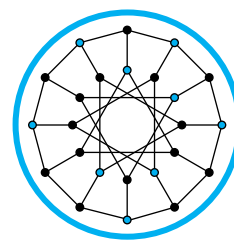
Clustering
algorithm
configuration



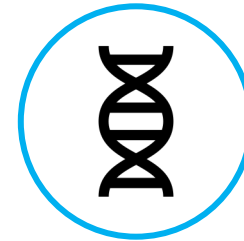
Integer programming
algorithm
configuration



Selling mechanism
configuration



Greedy
algorithm
configuration



Computational biology
algorithm
configuration



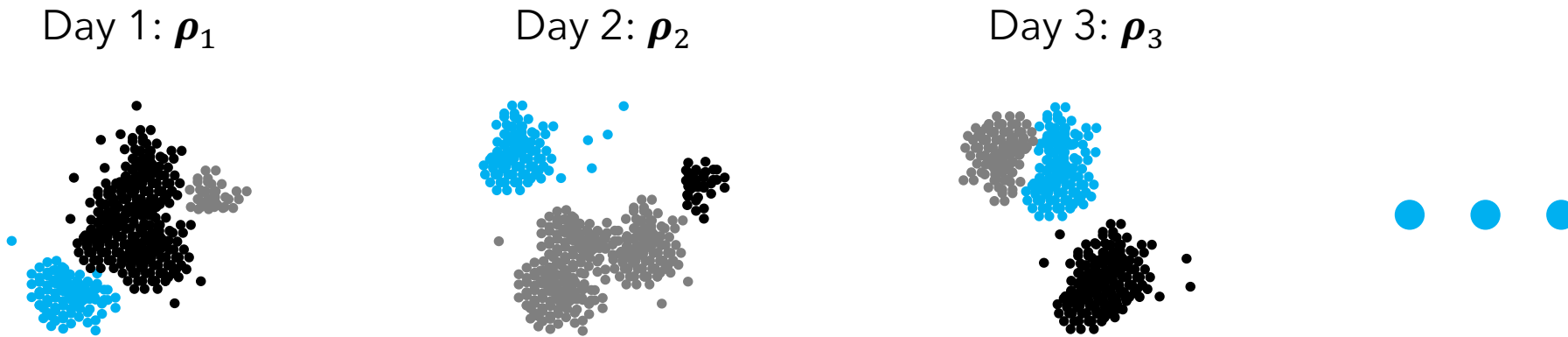
Voting mechanism
configuration

Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
 - i. Model
 - ii. Piecewise-structured algorithmic performance
 - iii. Main result
 - iv. Application: Sequence alignment
 - v. Online algorithm configuration**
2. Algorithms with predictions

Online algorithm configuration

What if inputs are not i.i.d., but even adversarial?



Goal: Compete with best parameter setting in hindsight

- Impossible in the worst case
- Under what conditions is online configuration possible?

Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
- 2. Algorithms with predictions**

Algorithms with predictions

Assume you have some **predictions** about your problem, e.g.:



Probability any given element is in a huge database

Kraska et al., SIGMOD'18; Mitzenmacher, NeurIPS'18

In caching, the next time you'll see an element

Lykouris, Vassilvitskii, ICML'18

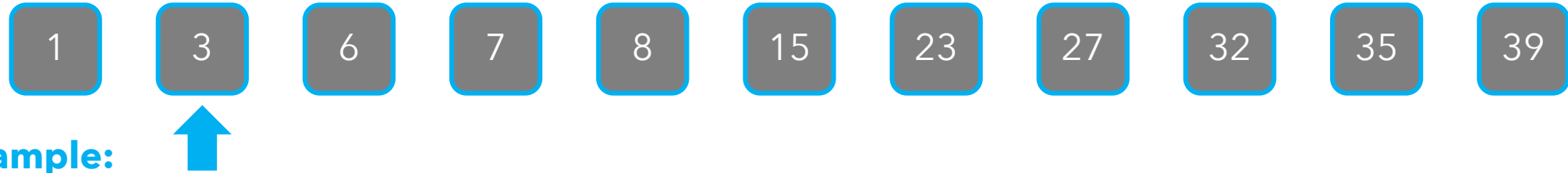
Main question:

How to use predictions to improve algorithmic performance?

Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
2. Algorithms with predictions
 - a. **Searching a sorted array**
 - b. Online algorithms
 - c. Additional research

Example: Searching in a sorted array

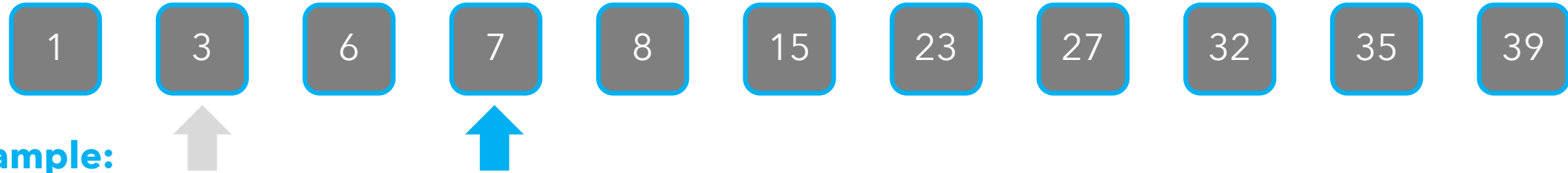


Example:

- $q = 8$
- $h(q) = 2$

- **Goal:** Given query q & sorted array A , find q 's index (if q in A)
- **Predictor:** $h(q) =$ guess of q 's index
- **Algorithm:** Check $A[h(q)]$. If q is there, return $h(q)$. Else:

Example: Searching in a sorted array

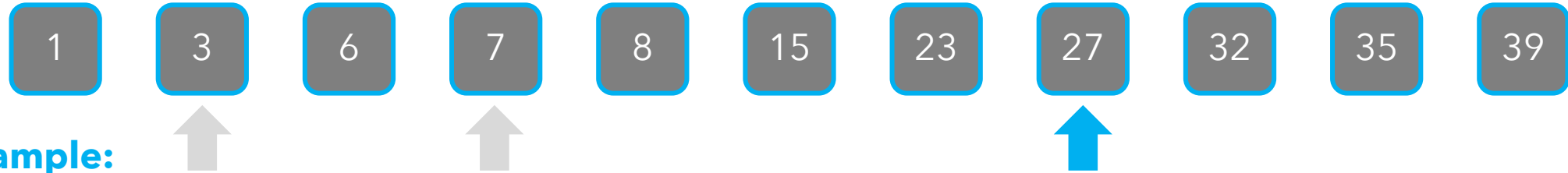


Example:

- $q = 8$
- $h(q) = 2$

- **Goal:** Given query q & sorted array A , find q 's index (if q in A)
- **Predictor:** $h(q) =$ guess of q 's index
- **Algorithm:** Check $A[h(q)]$. If q is there, return $h(q)$. Else:
 - If $q > A[h(q)]$, check $A[h(q) + 2^i]$ for $i > 1$ until find something larger

Example: Searching in a sorted array

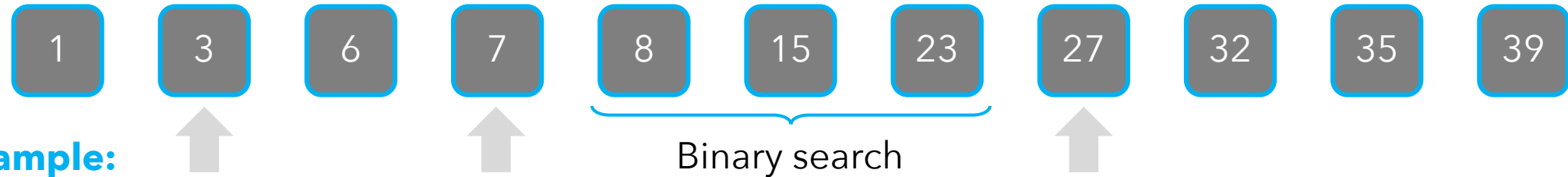


Example:

- $q = 8$
- $h(q) = 2$

- **Goal:** Given query q & sorted array A , find q 's index (if q in A)
- **Predictor:** $h(q) =$ guess of q 's index
- **Algorithm:** Check $A[h(q)]$. If q is there, return $h(q)$. Else:
 - If $q > A[h(q)]$, check $A[h(q) + 2^i]$ for $i > 1$ until find something larger
 - Do binary search on interval $(h(q) + 2^{i-1}, h(q) + 2^i)$

Example: Searching in a sorted array



Example:

- $q = 8$
- $h(q) = 2$

- **Goal:** Given query q & sorted array A , find q 's index (if q in A)
- **Predictor:** $h(q) =$ guess of q 's index
- **Algorithm:** Check $A[h(q)]$. If q is there, return $h(q)$. Else:
 - If $q > A[h(q)]$, check $A[h(q) + 2^i]$ for $i > 1$ until find something larger
 - Do binary search on interval $(h(q) + 2^{i-1}, h(q) + 2^i)$
 - If $q < A[h(q)]$, symmetric

Example: Searching in a sorted array



Analysis:

- Let $t(q)$ be index of q in A or of smallest element larger than q
- Runtime is $O(\log|t(q) - h(q)|)$:
 - Prediction error
 - Finding larger/smaller element takes $O(\log|t(q) - h(q)|)$ steps
 - Binary search takes $O(\log|t(q) - h(q)|)$ steps
- Better predictions lead to **better runtime**
- Runtime **never worse than worst-case** $O(\log|A|)$

Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
2. Algorithms with predictions
 - a. Searching a sorted array
 - b. Online algorithms**
 - c. Additional research

Purohit, Svitkina, Kumar, NeurIPS'18

Online algorithms

Full input not revealed upfront, but at some later stage, e.g.:

Matching: nodes of a graph arrive over time

Must irrevocably decide whether to match a node when it arrives

Caching: memory access requests arrive over time

Must decide what to keep in cache

Scheduling: job lengths not revealed until they terminate

Must decide which jobs to schedule when

Competitive ratio (CR)

Standard measure of online algorithm's performance:

$$\text{CR} = \frac{\text{ALG}}{\text{OPT}}$$

Offline optimal solution that knows the entire input

E.g., in matching:

$$\text{CR} = \frac{\text{weight of algorithm's matching}}{\text{maximum weight matching}}$$

Online algorithms

Full input not revealed upfront, but at some later stage

What if algorithm receives some **predictions** about input?

- **Online advertising**

e.g., Mahdian et al. [EC'07]; Devanur, Hayes [EC'09]; Muñoz Medina, Vassilvitskii [NeurIPS'17]

- **Caching**

e.g., Lykouris, Vassilvitskii [ICML'18]

- **Data structures**

e.g., Mitzenmacher [NeurIPS'18]

- ...

Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
2. Algorithms with predictions
 - a. Searching a sorted array
 - b. Online algorithms
 - i. Overview
 - ii. Ski rental problem**
 - iii. Job scheduling
 - c. Additional research

Ski rental problem

Family of problems that revolve around a decision:

- Incur a **recurring expense**, or
- Pay a **single fee** that eliminates the ongoing cost



Ski rental problem

Problem: Skier will ski for unknown number of days

- Can either **rent each day** for \$1/day or **buy** for \$ b
- E.g., if ski for 5 days and then buy, total price is $5 + b$

If ski x days, **optimal clairvoyant** strategy pays $\text{OPT} = \min\{x, b\}$

Breakeven strategy: Rent for $b - 1$ days, then buy

- $\text{CR} = \frac{\text{ALG}}{\text{OPT}} = \frac{x\mathbf{1}_{\{x < b\}} + (b-1+b)\mathbf{1}_{\{x \geq b\}}}{\min\{x, b\}} < 2$ (best deterministic)
- Randomized alg. $\text{CR} = \frac{e}{e-1}$ [Karlin et al., Algorithmica '94]



Ski rental problem

Prediction y of number of skiing days, error $\eta = |x - y|$

Baseline: Buy at beginning if $y > b$, else rent all days

Theorem: $\text{ALG} \leq \text{OPT} + \eta$

If y small but $x \gg b$, CR can be unbounded



Deterministic algorithm

Prediction y of number of skiing days, error $\eta = |x - y|$

Algorithm (with parameter $\lambda \in [0,1]$):

If $y \geq b$, buy on start of day $\lceil \lambda b \rceil$; else buy on start of day $\left\lceil \frac{b}{\lambda} \right\rceil$

- If **really trust** predictions: set $\lambda = 0$
Equivalent to blindly following predictions
- If **don't trust** predictions: set $\lambda = 1$
Equivalent to running the worst-case algorithm

Deterministic algorithm

Prediction y of number of skiing days, error $\eta = |x - y|$

Algorithm (with parameter $\lambda \in [0,1]$):

If $y \geq b$, buy on start of day $\lceil \lambda b \rceil$; else buy on start of day $\lceil \frac{b}{\lambda} \rceil$

Theorem: Algorithm has $\text{CR} \leq \min \left\{ \frac{1+\lambda}{\lambda}, 1 + \lambda + \frac{\eta}{(1-\lambda)\text{OPT}} \right\}$

- If predictor is perfect ($\eta = 0$), **CR is small** ($\leq 1 + \lambda$)
- No matter how big η is, setting $\lambda = 1$ **recovers baseline** $\text{CR} = 2$

Deterministic algorithm

Theorem: Algorithm has $CR \leq \min \left\{ \frac{1+\lambda}{\lambda}, 1 + \lambda + \frac{\eta}{(1-\lambda)OPT} \right\}$

Proof sketch: If $y \geq b$, buys on start of day $[\lambda b]$

$$\frac{ALG}{OPT} = \begin{cases} \frac{x}{x} & \text{if } x < [\lambda b] \end{cases}$$

Deterministic algorithm

Theorem: Algorithm has $\text{CR} \leq \min \left\{ \frac{1+\lambda}{\lambda}, 1 + \lambda + \frac{\eta}{(1-\lambda)\text{OPT}} \right\}$

Proof sketch: If $y \geq b$, buys on start of day $\lceil \lambda b \rceil$

$$\frac{\text{ALG}}{\text{OPT}} = \begin{cases} \frac{x}{x} & \text{if } x < \lceil \lambda b \rceil \\ \frac{\lceil \lambda b \rceil - 1 + b}{x} & \text{if } \lceil \lambda b \rceil \leq x \leq b \end{cases}$$

Deterministic algorithm

Theorem: Algorithm has $CR \leq \min \left\{ \frac{1+\lambda}{\lambda}, 1 + \lambda + \frac{\eta}{(1-\lambda)OPT} \right\}$

Proof sketch: If $y \geq b$, buys on start of day $\lceil \lambda b \rceil$

$$\frac{ALG}{OPT} = \begin{cases} \frac{x}{x} & \text{if } x < \lceil \lambda b \rceil \\ \frac{\lceil \lambda b \rceil - 1 + b}{x} & \text{if } \lceil \lambda b \rceil \leq x \leq b \\ \frac{\lceil \lambda b \rceil - 1 + b}{b} & \text{if } x \geq b \end{cases}$$

Worst when $x = \lceil \lambda b \rceil$ and $CR = \frac{b + \lceil \lambda b \rceil - 1}{\lceil \lambda b \rceil} \leq \frac{1+\lambda}{\lambda}$; similarly for $y < b$

Design principals

Consistency:

- Predictions are perfect \Rightarrow recover offline optimal
- Algorithm is α -consistent if $CR \rightarrow \alpha$ as error $\eta \rightarrow 0$

Robustness:

- Predictions are terrible \Rightarrow no worse than worst-case
- Algorithm is β -consistent if $CR \leq \beta$ for all η

E.g., ski rental: $CR \leq \min \left\{ \frac{1+\lambda}{\lambda}, 1 + \lambda + \frac{\eta}{(1-\lambda)OPT} \right\}$

$(1 + \lambda)$ -consistent, $\left(\frac{1+\lambda}{\lambda}\right)$ -robust

Bounds are tight [Gollapudi, Panigrahi, ICML'19; Angelopoulos et al., ITCS'20]



Randomized algorithm

if $y \geq b$:

Let $k \leftarrow \lfloor \lambda b \rfloor$

For $i \in [k]$, define $q_i \leftarrow \left(\frac{b-1}{b}\right)^{k-i} \frac{1}{b(1-(1-1/b)^k)}$

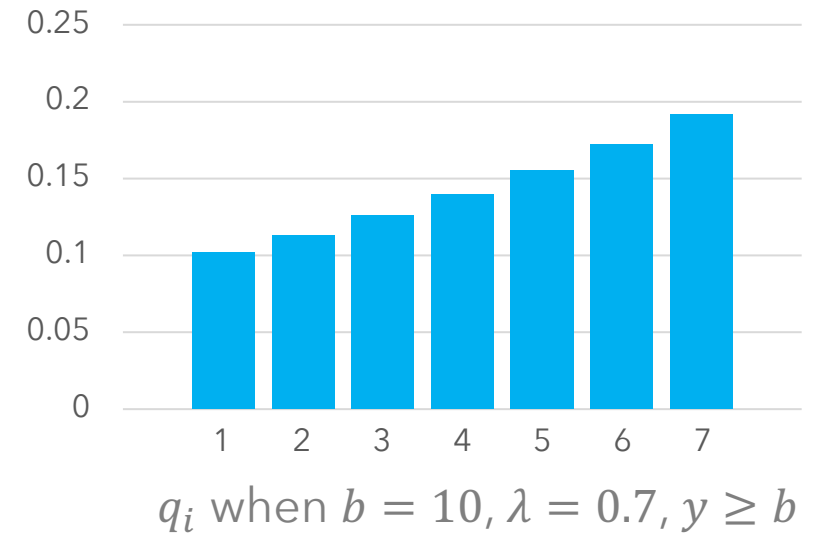
Buy on day $j \in [k]$ sampled from distribution defined by q_1, \dots, q_k

else

Let $\ell \leftarrow \lfloor \frac{b}{\lambda} \rfloor$

For $i \in [k]$, define $q_i \leftarrow \left(\frac{b-1}{b}\right)^{\ell-i} \frac{1}{b(1-(1-1/b)^\ell)}$

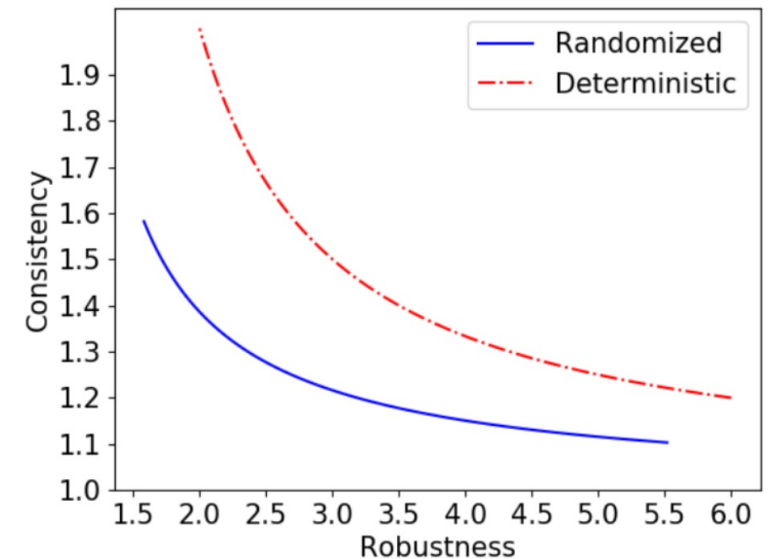
Buy on day $j \in [\ell]$ sampled from distribution defined by q_1, \dots, q_ℓ



Randomized algorithm

Theorem: $CR \leq \min \left\{ \frac{1}{1 - \exp(-(\lambda - 1/b))}, \frac{\lambda}{1 - \exp(-\lambda)} \left(1 + \frac{\eta}{OPT} \right) \right\}$

- $\left(\frac{\lambda}{1 - \exp(-\lambda)} \right)$ -consistent, $\left(\frac{1}{1 - \exp(-(\lambda - 1/b))} \right)$ -robust
- Bounds are **tight** [Wei, Zhang, NeurIPS'20]

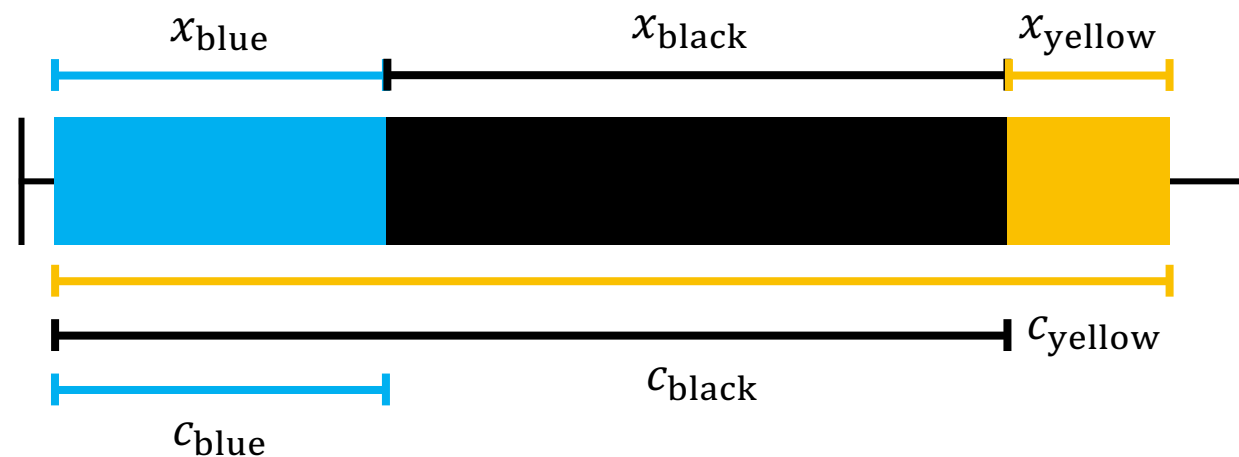


Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
2. Algorithms with predictions
 - a. Searching a sorted array
 - b. Online algorithms
 - i. Overview
 - ii. Ski rental problem
 - iii. Job scheduling**
 - c. Additional research

Job scheduling

- Task: schedule n jobs on a single machine
- Job j has **unknown** processing time x_j
- Goal: minimize **sum of completion times** of the jobs
i.e., if job j completes at time c_j , goal is to minimize $\sum c_j$



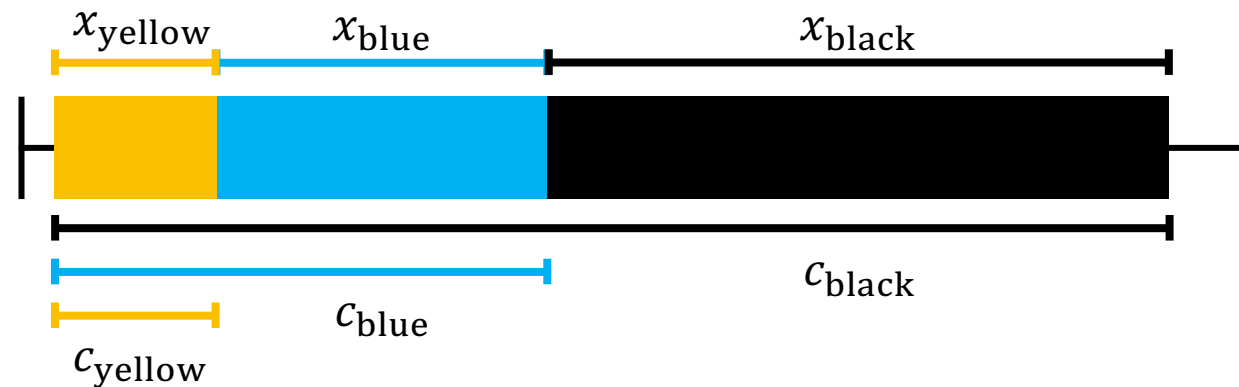
Job scheduling

- Task: schedule n jobs on a single machine
- Job j has **unknown** processing time x_j
- Goal: minimize **sum of completion times** of the jobs
i.e., if job j completes at time c_j , goal is to minimize $\sum c_j$
- Can switch between jobs



Job scheduling

Optimal solution if processing times x_j 's are known:
schedule jobs in increasing order of x_j



If $x_1 \leq \dots \leq x_n$,

$$\text{OPT} = \sum_{i=1}^n \sum_{j=1}^i x_j$$

Round robin

Algorithm with a competitive ratio of 2: **round robin**

Schedule 1 unit of time per remaining job, round-robin



Round-robin over k jobs \equiv run jobs simultaneously at rate of $\frac{1}{k}$



Algorithms-with-predictions approach

Predictions y_1, \dots, y_n of x_1, \dots, x_n with $\eta = \sum_{i=1}^n |y_i - x_i|$

If **really trust** predictions: schedule in increasing order of y_i
"Shortest predicted job first (SPJF)"

If **don't trust** predictions: round-robin (RR)

Algorithms-with-predictions approach

Algorithm: Preferential round-robin (with parameter $\lambda \in (0,1)$)

Run SPJF and RR **simultaneously**

- SPJF at a rate λ
- RR at a rate $1 - \lambda$

Example: $\lambda = \frac{1}{2}$, 3 jobs, shortest predicted job is **blue** job



- **Blue** job at a rate of $\lambda + (1 - \lambda) \cdot \frac{1}{3} = \frac{2}{3}$
- **Yellow** job at a rate of $(1 - \lambda) \cdot \frac{1}{3} = \frac{1}{6}$
- **Black** job at a rate of $(1 - \lambda) \cdot \frac{1}{3} = \frac{1}{6}$

Preferential round-robin

Algorithm: Preferential round-robin (with parameter $\lambda \in (0,1)$)

Run SPJF and RR **simultaneously**

- SPJF at a rate λ
- RR at a rate $1 - \lambda$

Theorem: Algorithm's competitive ratio is

$$\text{CR} \leq \min \left\{ \underbrace{\frac{1}{\lambda} \left(1 + \frac{2\eta}{n} \right)}_{\text{CR of SPJF}}, \underbrace{\frac{1}{1-\lambda} \cdot 2}_{\text{CR of RR}} \right\}$$

So it's $\frac{1}{\lambda}$ -consistent, $\frac{2}{1-\lambda}$ -robust

Outline (theoretical guarantees)

1. Statistical guarantees for algorithm configuration
2. Algorithms with predictions
 - a. Searching a sorted array
 - b. Online algorithms
 - c. Additional research**

Just scratched the surface

Online advertising

Mahdian, Nazerzadeh, Saberi, EC'07;
Devanur, Hayes, EC'09; Medina,
Vassilvitskii, NeurIPS'17; ...

Caching

Lykouris, Vassilvitskii, ICML'18; Rohatgi,
SODA'19; Wei, APPROX-RANDOM'20; ...

Frequency estimation

Hsu, Indyk, Katabi, Vakilian, ICLR'19; ...

Learning low-rank approximations

Indyk, Vakilian, Yuan, NeurIPS'19; ...

Scheduling

Mitzenmacher, ITCS'20; Moseley,
Vassilvitskii, Lattanzi, Lavastida, SODA'20; ...

[algorithms-with-predictions.github.io](https://github.com/ellenkit/algorithm-with-predictions)

Closely related: the “predict-then-optimize” framework

Elmachtoub, Grigas, Management Science '22; Elmachtoub et al., ICML'20; ...

Summary

1 Applied techniques

- a. Graph neural networks
 - a. Neural algorithmic alignment
 - b. Variable selection for integer programming
- b. Learning greedy heuristics with RL

2 Theoretical guarantees

- a. Statistical guarantees for algorithm configuration
- b. Algorithms with predictions

3 Future directions

Outline (future directions)

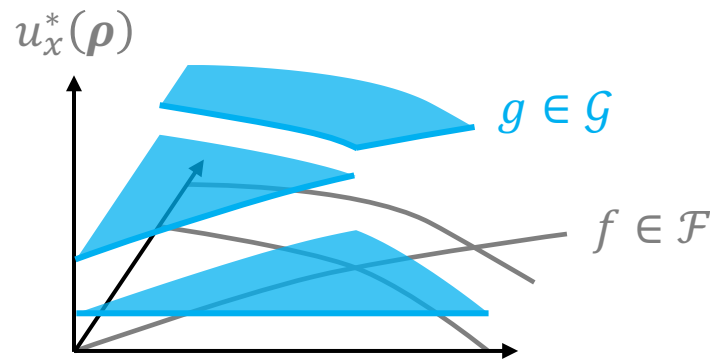
- 1. Tighter statistical bounds**
2. Multi-task algorithm design: Knowledge transfer
3. Size generalization
4. ML as a toolkit for theory

Future work: Tighter statistical bounds

WHP $\forall \rho$, **avg** utility over training set - **exp** utility $\leq \epsilon$

given training set of size $\tilde{O}\left(\frac{1}{\epsilon^2} (\text{Pdim}(\mathcal{G}^*) + \text{VCdim}(\mathcal{F}^*) \log k)\right)$

Number of boundary functions



k is often exponential
Can lead to large bounds

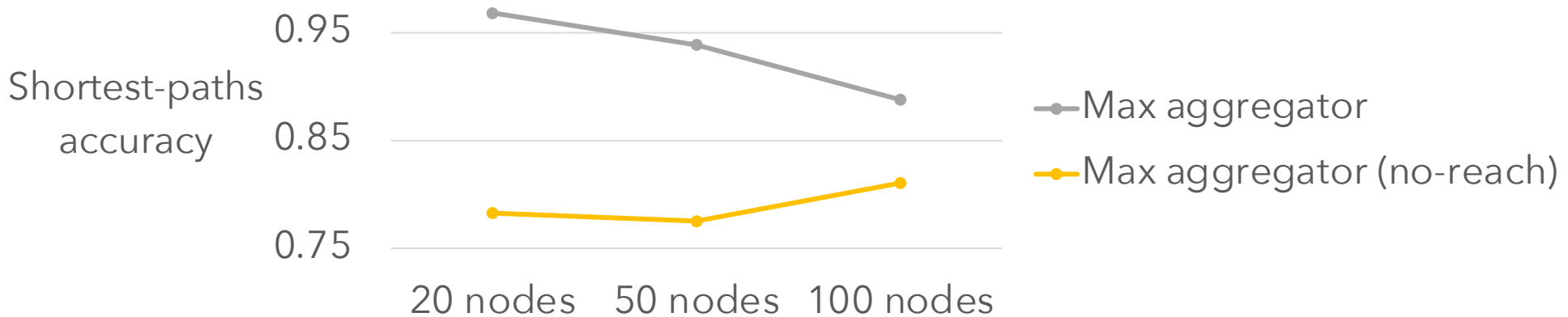
I expect this can sometimes be avoided!
Would require more information about duals

Outline (future directions)

1. Tighter statistical bounds
- 2. Multi-task algorithm design: Knowledge transfer**
3. Size generalization
4. ML as a toolkit for theory

Future work: Knowledge transfer

- Training a GNN to solve multiple related problems... can sometimes lead to better **single-task** performance
- E.g., training reachability and shortest-paths (grey line) v.s. just training shortest-paths (**yellow line**)



Future work: Knowledge transfer

- Training a GNN to solve multiple related problems...
can sometimes lead to better **single-task** performance
- Can we understand **theoretically** why this happens?
 - For which sets of algorithms can we expect **knowledge transfer**?

Outline (future directions)

1. Tighter statistical bounds
2. Multi-task algorithm design: Knowledge transfer
- 3. Size generalization**
4. ML as a toolkit for theory

Future work: Size generalization

Machine-learned algorithms can **scale to larger instances**

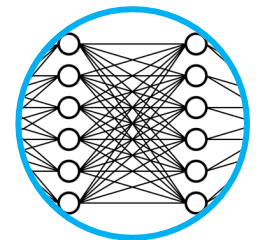
Applied research: Dai et al., NeurIPS'17; Veličković, et al., ICLR'20; ...

Goal: eventually, solve problems **no one's ever been able to solve**

However, size generalization is not immediate! It depends on:

- The **machine-learned algorithm**

Is the algorithm scale sensitive?



Example [Xu et al., ICLR'21]:

- Algorithms represented by GNNs **do generalize**
- Algs represented by MLPs **don't generalize** across size

Future work: Size generalization

Machine-learned algorithms can **scale to larger instances**

Applied research: Dai et al., NeurIPS'17; Veličković, et al., ICLR'20; ...

Goal: eventually, solve problems **no one's ever been able to solve**

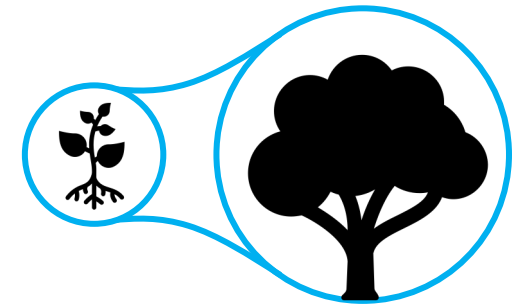
However, size generalization is not immediate! It depends on:

- The **machine-learned algorithm**

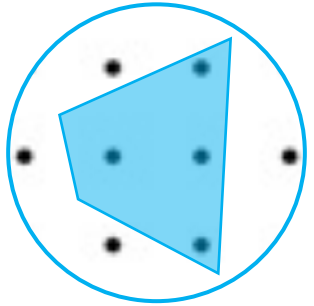
Is the algorithm scale sensitive?

- The **problem instances**

As size scales, what features must be preserved?



Future work: Size generalization



Can you:

1. **Shrink** a set of big integer programs
graphs

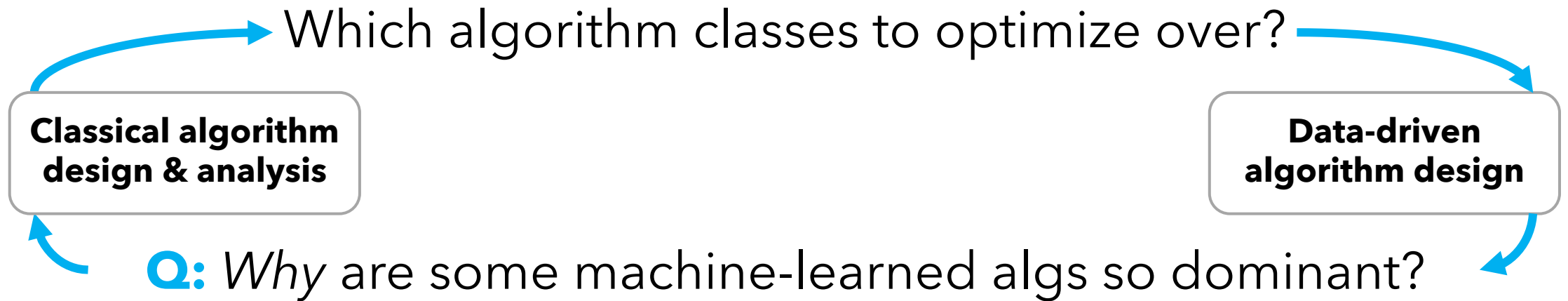
...

2. **Learn** a good algorithm on the **small** instances
3. **Apply** what you learned to the **big** instances?

Outline (future directions)

1. Tighter statistical bounds
2. Multi-task algorithm design: Knowledge transfer
3. Size generalization
- 4. ML as a toolkit for theory**

Future work: ML as a toolkit for theory



E.g., Dai et al. [NeurIPS'17] write that their RL alg discovered:
"New and interesting" greedy strategies for MAXCUT and MVC
"which **intuitively make sense** but have **not been analyzed** before,"
thus could be a "good **assistive tool** for discovering new algorithms."

Thank you! Questions?