# Learning Probabilistic and Causal Models with(out) Imperfect Advice

PhD Defense

13 January 2025

Davin Choo

National University of Singapore

# How do we find words in a dictionary?
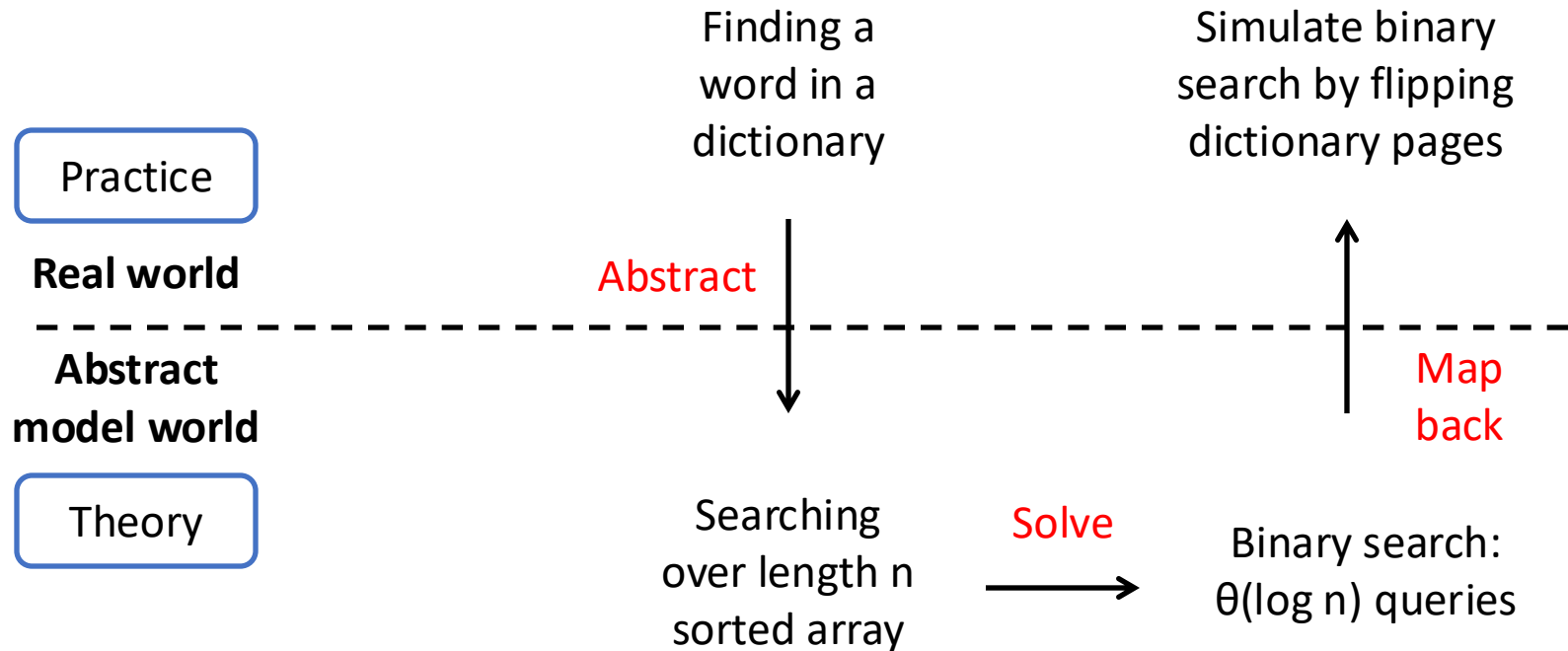
# How do we find words in a dictionary?



Linear search
O(n) pages

Binary search
O(log n) pages

imgflip.com

# A general problem-solving framework

# A general problem-solving framework

# A general problem-solving framework

- Complex setting
- Many nuances
- Possibly unseen problem

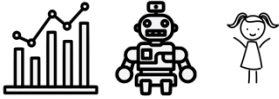**Real world**

- - - - - - - - - - - - - - - - -

**Abstract**
**model world**

- Simplified setting
- Generic problem framing
- Many plug-and-play solution concepts

<u>Two useful scientific models</u>

1) **Probabilistic models** for predictive tasks

2) **Causal models** for understanding

interventional effects on systems

# Side-information about problem instances



Real world

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

Abstract
model world

Finding a word in a dictionary
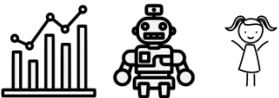
Simulate binary search by flipping dictionary pages

Abstract

Map back

Searching over length n sorted array
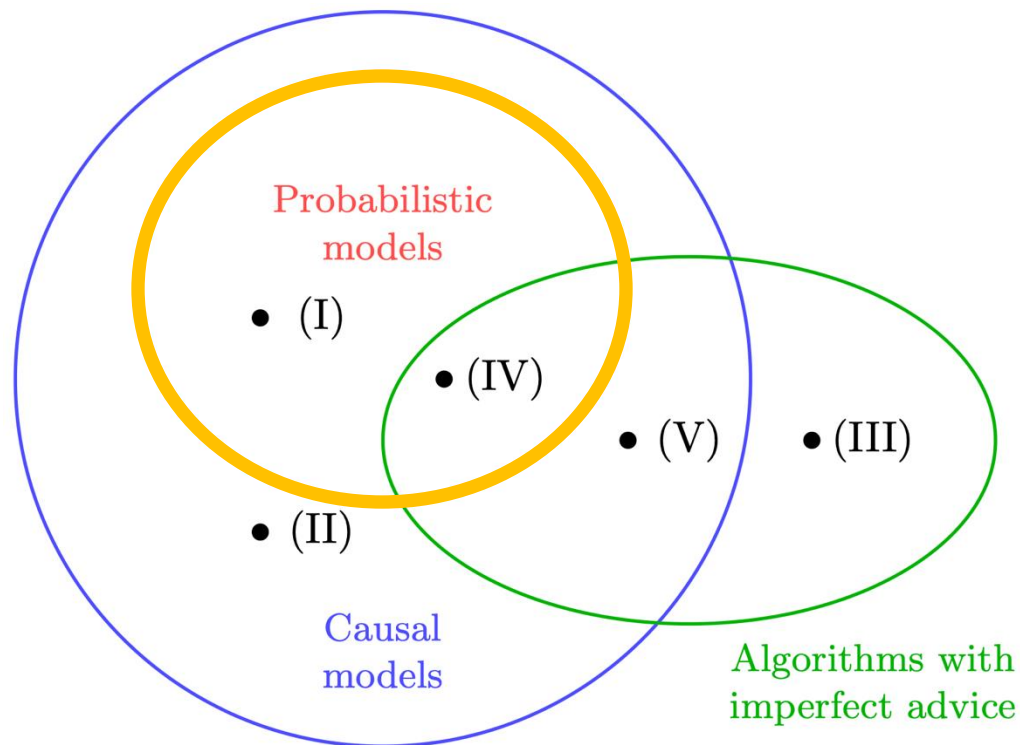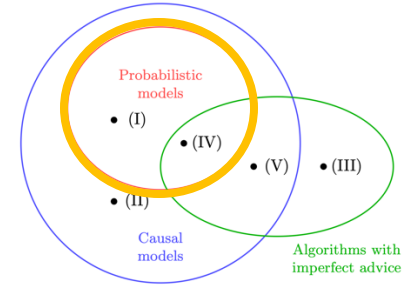
Solve

Binary search: θ(log n) queries

# Side-information about problem instances

# Main themes explored in my PhD thesis

# Main themes explored in my PhD thesis

# (I): Probabilistic models

- Classic results in statistics show asymptotic convergence of estimators in the limit of large data

- Probably Approximately Correct (PAC) learning model [Val84]
  - Given sample access to some underlying distribution $\mathcal{P}$, produce $\hat{\mathcal{P}}$ such that $\text{TV}(\mathcal{P}, \hat{\mathcal{P}}) \leq \varepsilon$ with probability $\geq 1 - \delta$

Probability mass, i.e. area under curve sums to 1

P

Q

Domain $\mathcal{X}$

$\text{TV}(P, Q)$

[Val84] Leslie G Valiant. *A theory of the learnable*. Communications of the ACM, 1984.
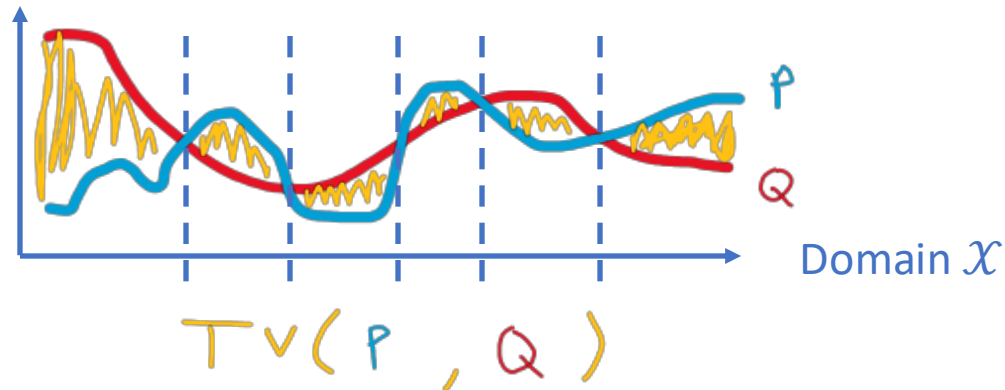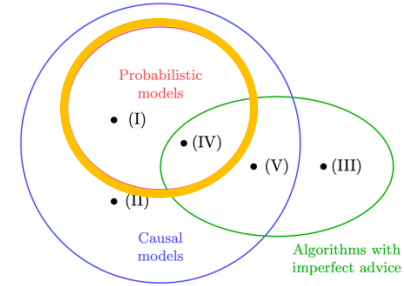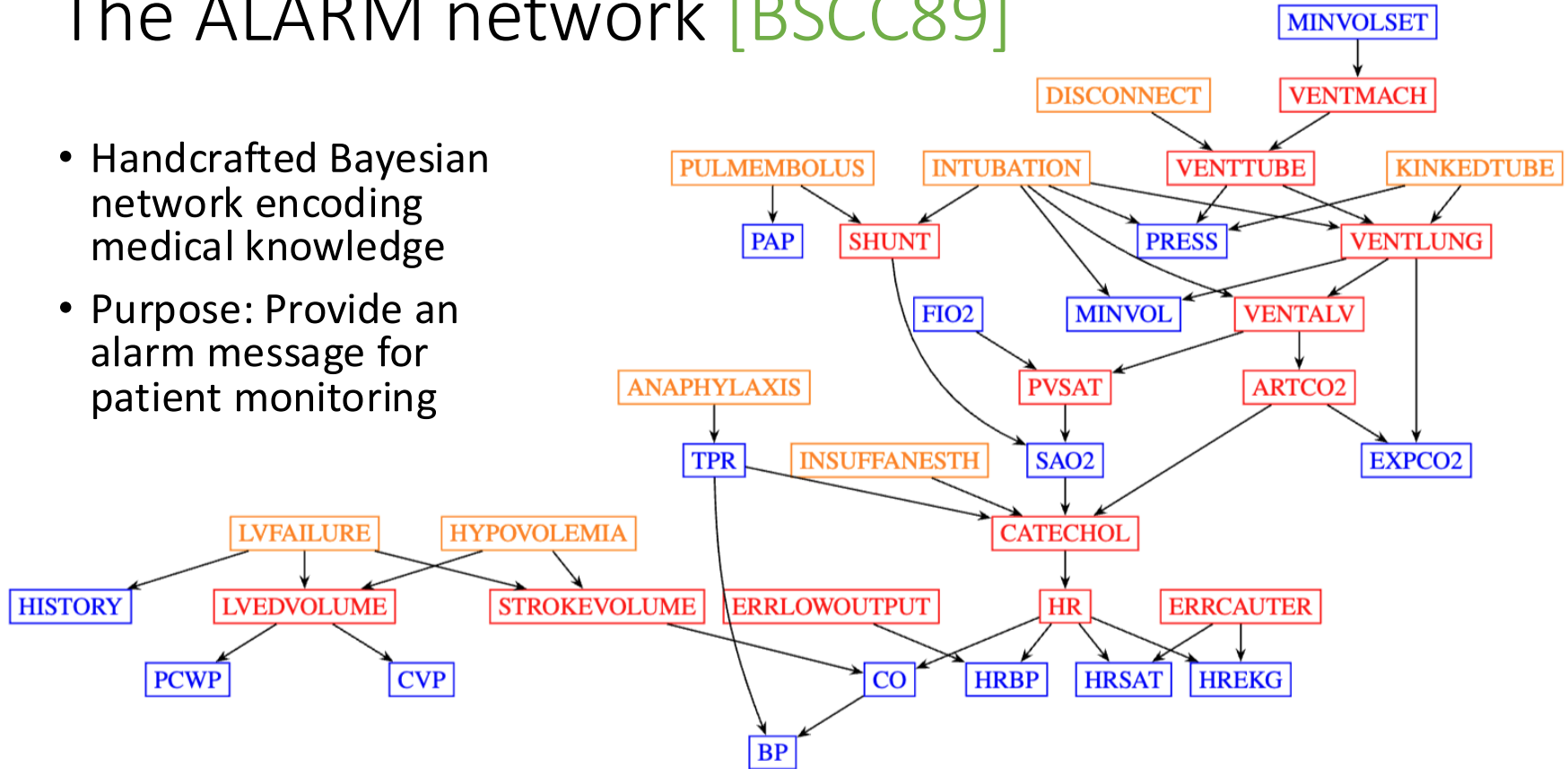
# (I): Probabilistic models



- Classic results in statistics show asymptotic convergence of estimators in the limit of large data

- Probably Approximately Correct (PAC) learning model [Val84]
  - Given sample access to some underlying distribution $\mathcal{P}$, produce $\hat{\mathcal{P}}$ such that $\text{TV}(\mathcal{P}, \hat{\mathcal{P}}) \leq \varepsilon$ with probability $\geq 1 - \delta$

- Bayesian networks [Pea88]
  - Probabilistic graphical model commonly used to model beliefs
  - 2 parts: graph + conditional distributions for each vertex
  - $\approx 2^{n^2}$ candidate directed acyclic graphs (DAGs), one of which is $\mathcal{G}^*$

[Pea88] Judea Pearl. *Probabilistic reasoning in intelligent systems*. Morgan Kaufmann, 1988.

# The ALARM network [BSCC89]

- Handcrafted Bayesian network encoding medical knowledge

- Purpose: Provide an alarm message for patient monitoring



[BSCC89] Ingo A Beinlich, Henri Jacques Suermondt, R Martin Chavez, and Gregory F Cooper. *The alarm monitoring system: A case study with two probabilistic inference techniques for belief network*. Second European Conference on Artificial Intelligence in Medicine (AIME), 1989

# The ALARM network [BSCC89]

**A sample consultation**

ALARM is a data-driven system. Simulating an anesthesia monitor, ALARM accepts a set of physiologic measurements. An example would be as follows: blood pressure 120/80 mmHg, heart rate 80/min, inspired oxygen concentration 50%, tidal volume 500 ml, respiratory rate 10/min, breathing pressure 50 mbar, and measured minute ventilation 1.2 l/min. These measurements are categorized into 'low', 'normal', 'high', etc. and text messages are generated when measurements are outside of their normal range. These messages will then appear in the *Warning* and *Caution* fields of the monitor depending on their importance *(Fig. 3)*. In the given example, the high breathing pressure of 50 mbar imposes a direct danger to the patient and a warning is issued. The low minute ventilation is less immediate and is displayed as a caution only.
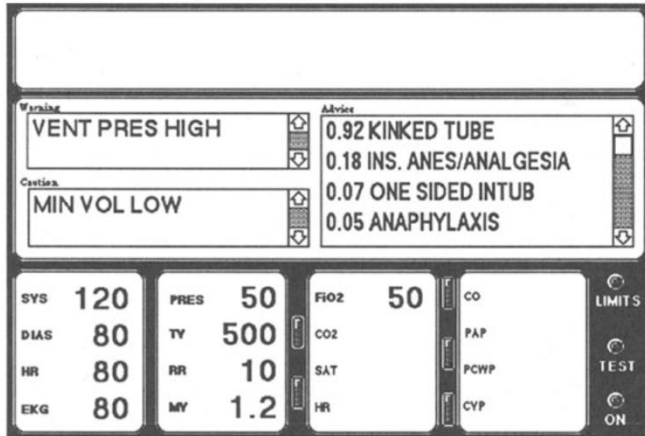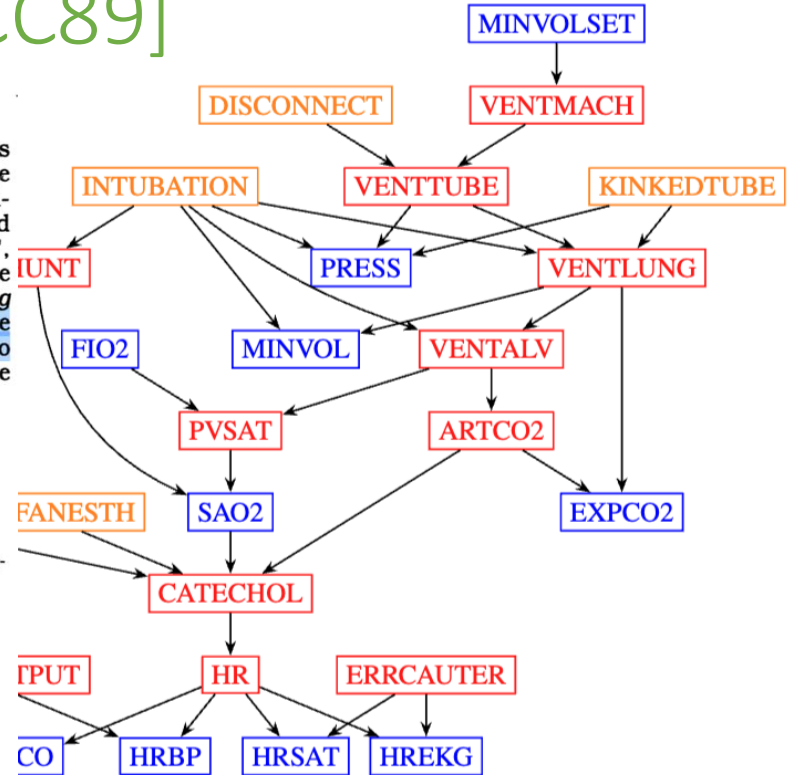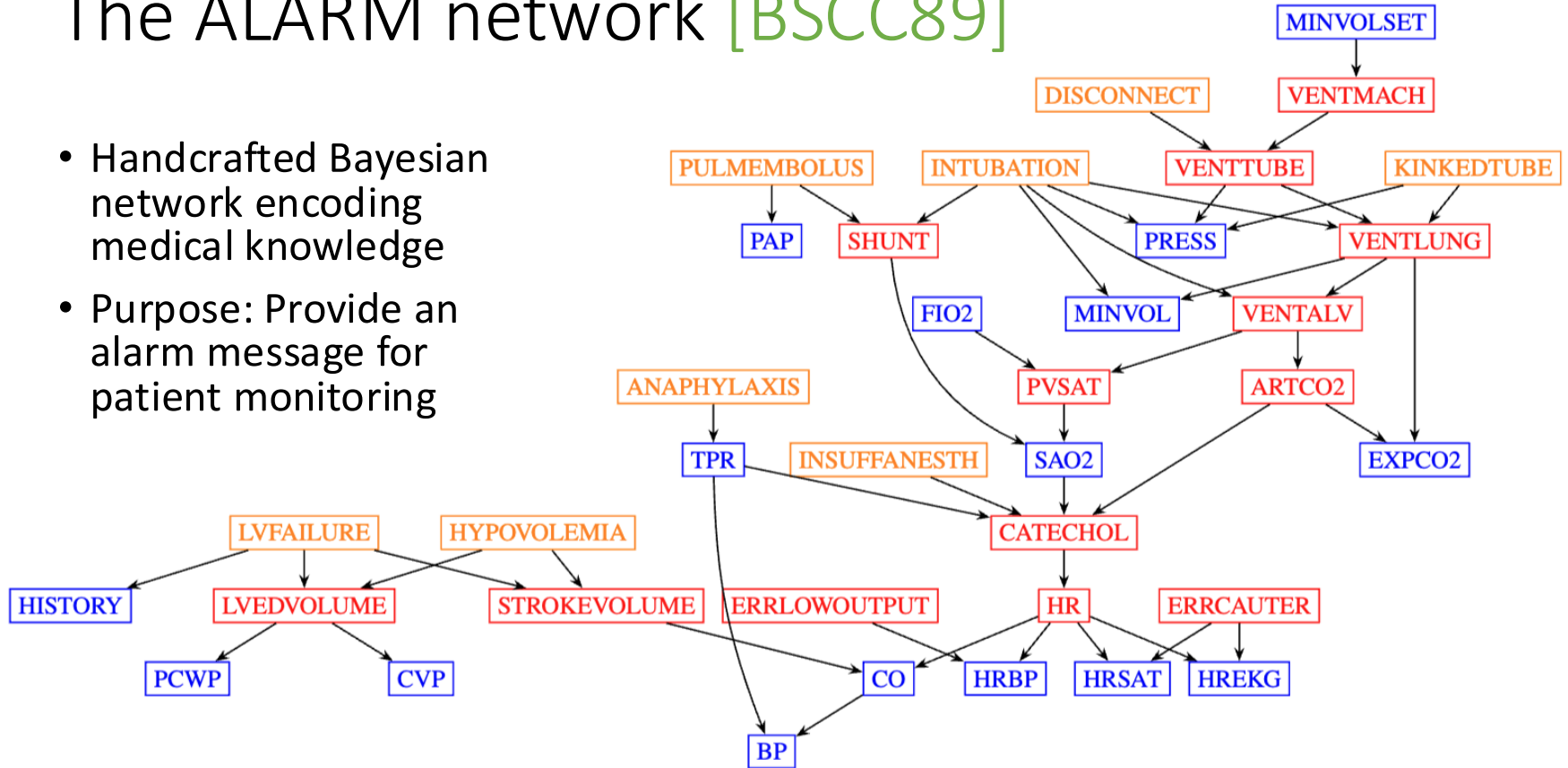
*Fig. 3*

ALARM simulates an anesthesia monitor. It takes patient measurements, displays warning and caution messages, and lists a differential diagnosis.
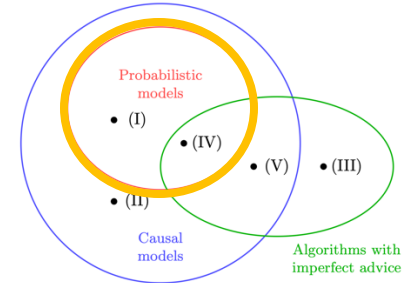
Warning
VENT PRES HIGH

Caution
MIN VOL LOW

Advice
0.92 KINKED TUBE
0.18 INS. ANES/ANALGESIA
0.07 ONE SIDED INTUB
0.05 ANAPHYLAXIS

| SYS | 120 | PRES | 50 | FiO2 | 50 | CO | | LIMITS |
| DIAS | 80 | TV | 500 | CO2 | | PAP | | |
| HR | 80 | RR | 10 | SAT | | PCWP | | TEST |
| EKG | 80 | MV | 1.2 | HR | | CYP | | ON |

[BSCC89] Ingo A Beinlich, Henri Jacques Suermondt, R Martin Chavez, and Gregory F Cooper. *The alarm monitoring system: A case study with two probabilistic inference techniques for belief network*. Second European Conference on Artificial Intelligence in Medicine (AIME), 1989

# The ALARM network [BSCC89]

- Handcrafted Bayesian network encoding medical knowledge

- Purpose: Provide an alarm message for patient monitoring



[BSCC89] Ingo A Beinlich, Henri Jacques Suermondt, R Martin Chavez, and Gregory F Cooper. *The alarm monitoring system: A case study with two probabilistic inference techniques for belief network*. Second European Conference on Artificial Intelligence in Medicine (AIME), 1989

# (I): Probabilistic models



- Suppose data distribution $\mathcal{P}$ is described by Bayesian network
  - NP-hard to find "score maximizing" DAG from data [Chi96] and to decide whether $\mathcal{P}$ can be described by a DAG with p parameters [CHM04]
  - Even under the promise that $\mathcal{P}$ can be described by a DAG with p parameters, it is NP-hard to find such a parameter-bounded DAG [B**C**GM25]
  - We also have some PAC-style finite sample results in learning the structure and parameters of Bayesian network for $\mathcal{P}$ [B**C**G+22, DDK**C**23, **C**YBC24]
  - Insight: If network's in-degree is bounded, we can use less samples

[Chi96] David Maxwell Chickering. *Learning Bayesian networks is NP-complete*. Lecture Notes in Statistics, vol 112, 1996
[CHM04] Max Chickering, David Heckerman, and Chris Meek. *Large-sample learning of Bayesian networks is NP-hard*. Journal of Machine Learning Research (JMLR), 2004
[B**C**GM25] Arnab Bhattacharyya, Davin Choo, Sutanu Gayen, Dimitrios Myrisiotis. *Learnability of Parameter-Bounded Bayes Nets*. AAAI Conference on Artificial Intelligence (AAAI), 2025
[B**C**G+22] Arnab Bhattacharyya, Davin Choo, Rishikesh Gajjala, Sutanu Gayen, Yuhao Wang. *Learning Sparse Fixed-Structure Gaussian Bayesian Networks*. International Conference on Artificial Intelligence and Statistics (AISTATS), 2022
[DDK**C**23] Yuval Dagan, Constantinos Daskalakis, Anthimos-Vardis Kandiros, Davin Choo. *Learning and Testing Latent-Tree Ising Models Efficiently*. Conference on Learning Theory (COLT), 2023
[**C**YBC24] Davin Choo, Joy Qiping Yang, Arnab Bhattacharyya, Clément L. Canonne. *Learning bounded degree polytrees with samples*. International Conference on Algorithmic Learning Theory (ALT), 2024
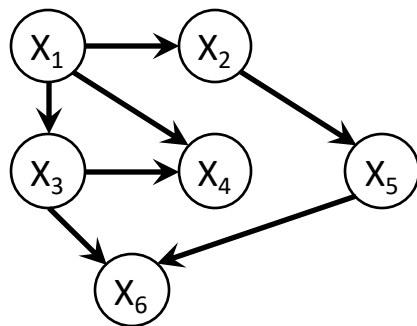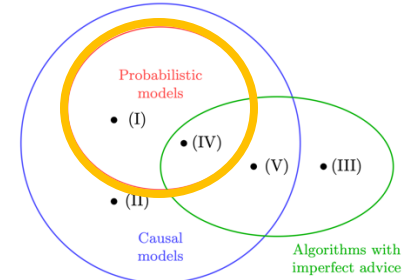
# A glimpse of [BCG+22]



Insight: If network's in-degree is bounded, we can use less samples

- Suppose we get i.i.d. samples from a linear DAG with Gaussian noise



$$\begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{pmatrix} = \begin{pmatrix} a_{1,1} & a_{1,2} & \ldots & a_{1,n} \\ a_{2,1} & a_{2,2} & 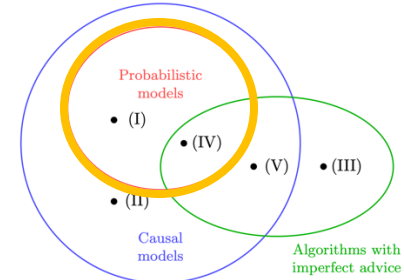\ldots & a_{2,n} \\ \vdots & \vdots & \ldots & \vdots \\ a_{n,1} & a_{n,2} & \ldots & a_{n,n} \end{pmatrix} \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{pmatrix} + \begin{pmatrix} \eta_1 \\ \eta_2 \\ \vdots \\ \eta_n \end{pmatrix}$$

$$X_i = \begin{cases} \eta_i + \sum_{X_j \in \mathrm{Pa}(X_i)} a_{i,j} X_j & \text{if } \mathrm{Pa}(X_i) \neq \emptyset \\ \eta_i & \text{if } \mathrm{Pa}(X_i) = \emptyset \end{cases}$$

# A glimpse of [B<u>C</u>G+22]



Insight: If network's in-degree is bounded, we can use less samples

- Suppose we get i.i.d. samples from a linear DAG with Gaussian noise



$$\begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{pmatrix} = \begin{pmatrix} a_{1,1} & a_{1,2} & \dots & a_{1,n} \\ a_{2,1} & a_{2,2} & \dots & a_{2,n} \\ \vdots & \vdots & \dots & \vdots \\ a_{n,1} & a_{n,2} & \dots & a_{n,n} \end{pmatrix} \begin{pmatrix} X_1 \\ X_2 \\ \vdots \\ X_n \end{pmatrix} + \begin{pmatrix} \eta_1 \\ \eta_2 \\ \vdots \\ \eta_n \end{pmatrix}$$
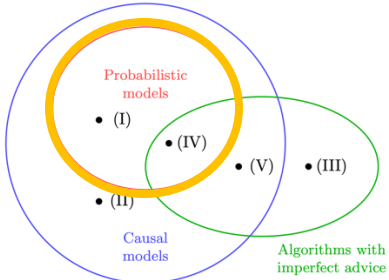
|  | $X_1$ | $X_2$ | $X_3$ | $X_4$ | $X_5$ | $X_6$ |
|---|---|---|---|---|---|---|
| Sample 1 | 1.24 | 1.08 | 0.229 | -0.846 | 0.307 | 1.201 |
| Sample 2 | -0.614 | 0.552 | 0.758 | 1.77 | 1.646 | 0.375 |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | ⋮ |

How many samples would we need to learn the coefficients and noise?

# A glimpse of [B**C**G+22]



Insight: If network's in-degree is bounded, we can use less samples

• Suppose we get i.i.d. samples from a linear DAG with Gaussian noise



$$X = AX + \eta$$
$$\Rightarrow X = (I_n - A)^{-1}\eta$$
$$\Rightarrow X \text{ is a multivariate Gaussian, in general}$$
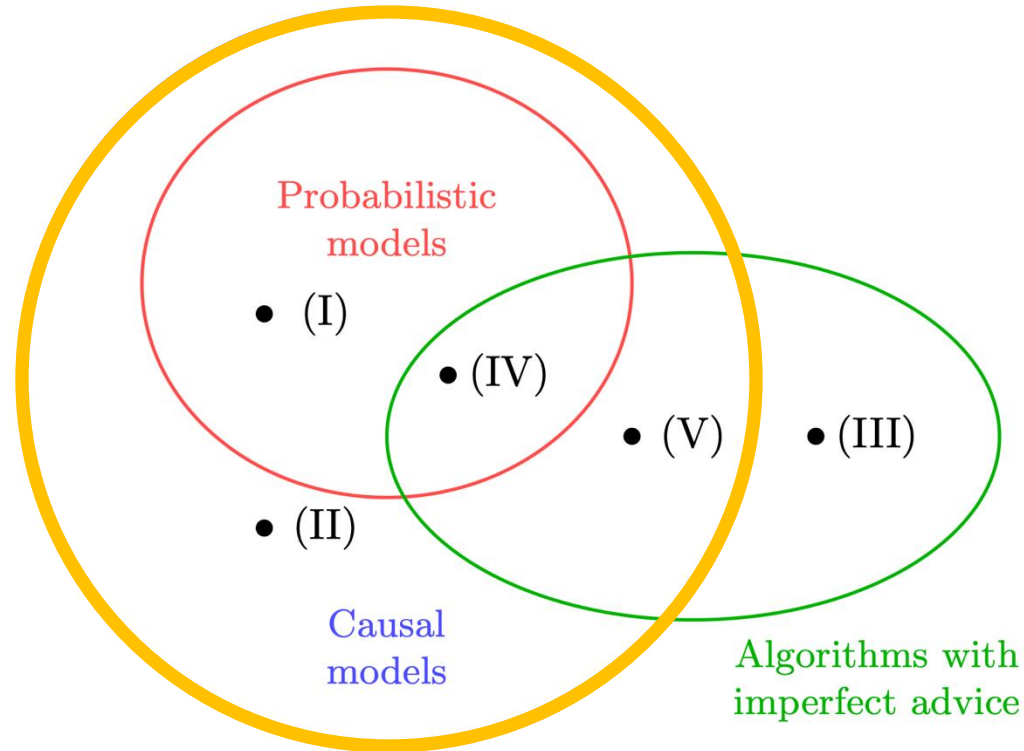$$\Rightarrow \text{Need } \widetilde{\Omega}\left(\frac{n^2}{\varepsilon^2}\right) \text{ i.i.d. samples to learn } X \text{ "}\varepsilon\text{-well"}$$

A rough intuition: All $n^2$ covariance matrix entries "matter", in general

# A glimpse of [BCG+22]



**Insight: If network's in-degree is bounded, we can use less samples**

- Turns out $\tilde{O}\left(\frac{nd}{\varepsilon^2}\right)$ samples suffice with just least squares at each node



$$\boldsymbol{X} = \boldsymbol{AX} + \boldsymbol{\eta}$$
$$\Rightarrow \boldsymbol{X} = (\boldsymbol{I}_n - \boldsymbol{A})^{-1}\boldsymbol{\eta}$$
$$\Rightarrow \boldsymbol{X} \text{ is a multivariate Gaussian, in general}$$

- Here, max in-degree $d = 2$   $\Rightarrow$ Need $\tilde{\Omega}\left(\frac{n^2}{\varepsilon^2}\right)$ i.i.d. samples to learn $\boldsymbol{X}$ "$\varepsilon$-well"

A rough intuition: All $n^2$ covariance matrix entries "matter", in general

# Main themes explored in my PhD thesis



Probabilistic models

• (I)

• (IV)

• (V)     • (III)

• (II)

Causal models

Algorithms with imperfect advice

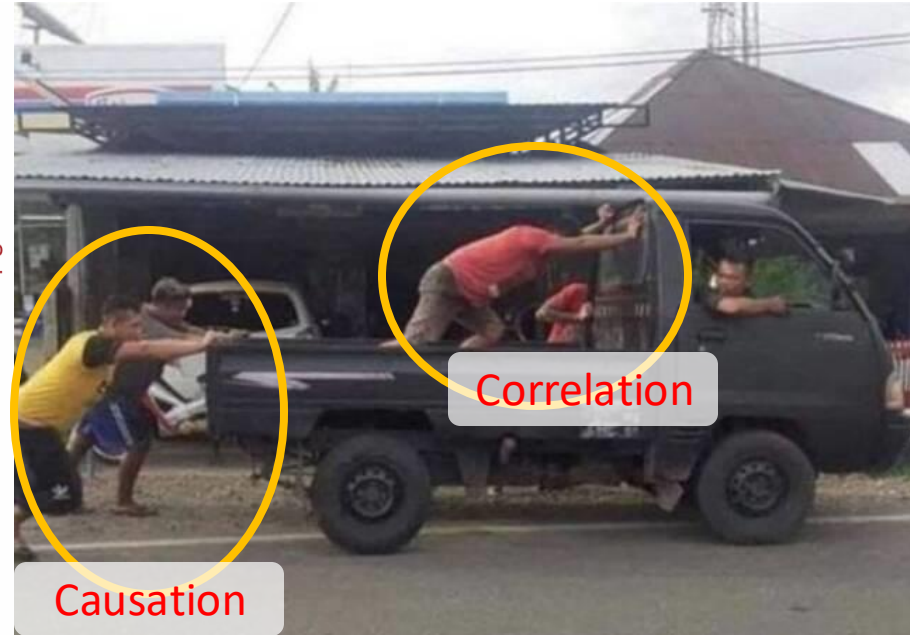# Correlation does not imply causation



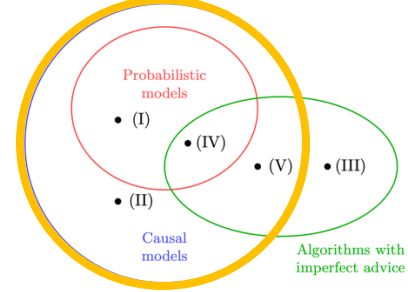**Robberies in Alaska**
correlates with
**Professor salaries in the US**

The robbery rate per 100,000 residents in Alaska · Source: FBI Criminal Justice Information Services

Average salary of full-time instructional faculty on 9-month contracts in degree-granting postsecondary institutions, by academic rank of Professor · Source: National Center for Education Statistics

2009-2021, r=0.922, r²=0.851, p<0.01 · tylervigen.com/spurious/correlation/2723



Correlation

Causation

# (II): Causal models



- Two fundamental problems in causal inference
  - Causal graph discovery: Recover true causal graph $\mathcal{G}^*$

  - Causal effect estimation: Estimate $\mathcal{P}(Y = y \mid do(X = x))$
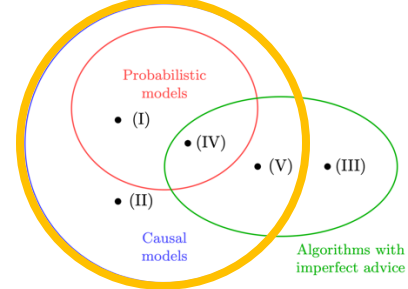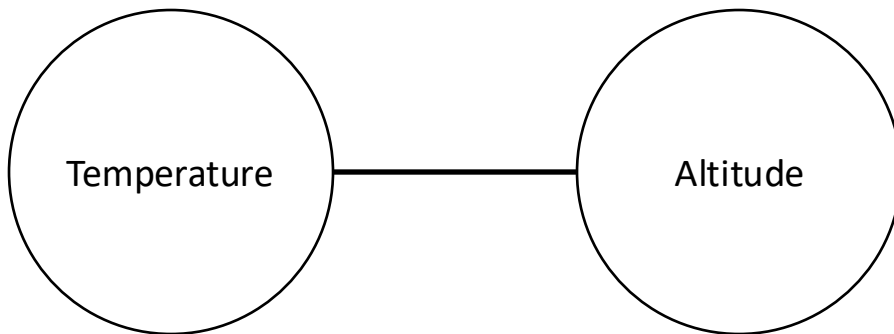
# (II): Causal models



- Two fundamental problems in causal inference
  - Causal graph discovery: Recover true causal graph $\mathcal{G}^*$
    - Even with infinite observational data, can only determine causal graph up to some equivalence class where all conditional independence relations agree
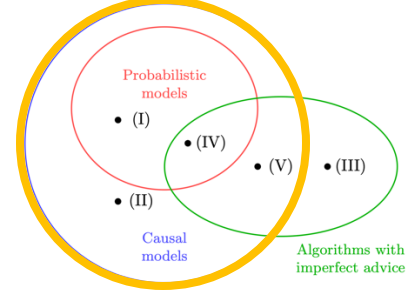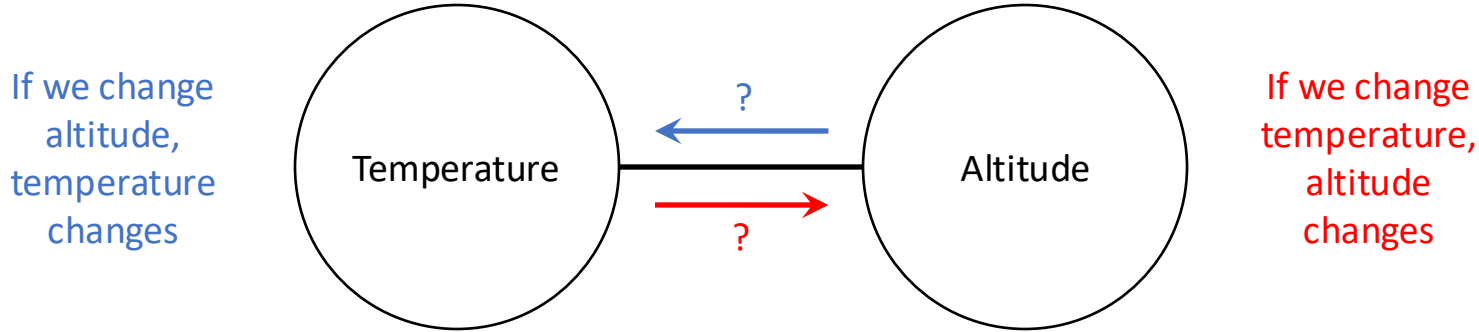    - Make distributional/structural assumptions or perform interventions/experiments!

# (II): Causal models



- Two fundamental problems in causal inference
  - Causal graph discovery: Recover true causal graph $\mathcal{G}^*$
    - Even with infinite observational data, can only determine causal graph up to some equivalence class where all conditional independence relations agree
    - Make distributional/structural assumptions or perform interventions/experiments!
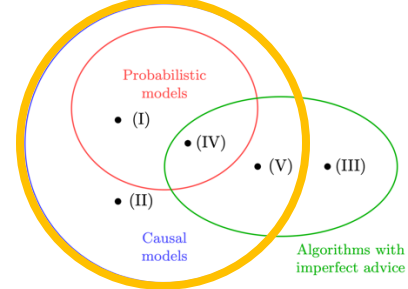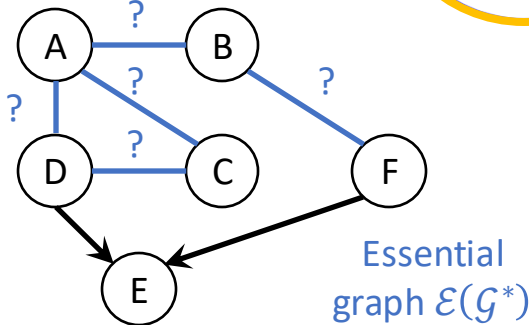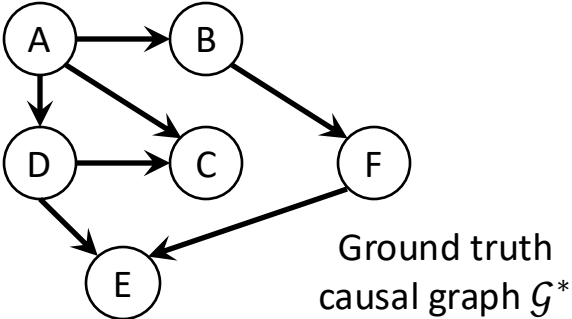
# (II): Causal models



- Two fundamental problems in causal inference
  - Causal graph discovery: Recover true causal graph $\mathcal{G}^*$
    - Even with infinite observational data, can only determine causal graph up to some equivalence class where all conditional independence relations agree
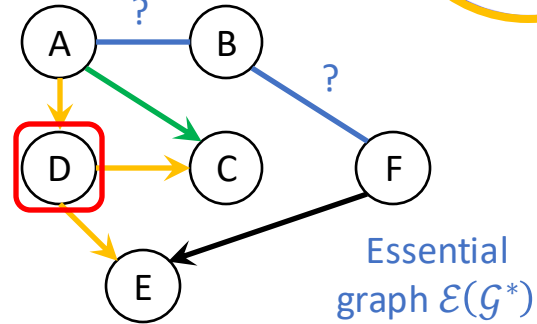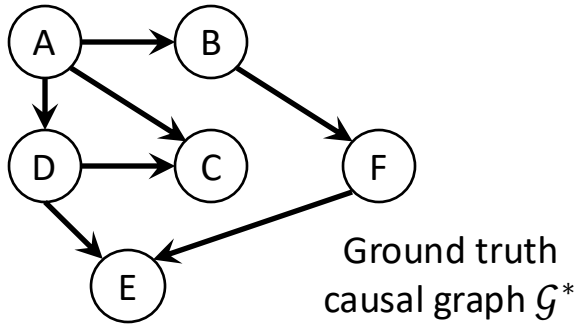    - Make distributional/structural assumptions or perform interventions/experiments!
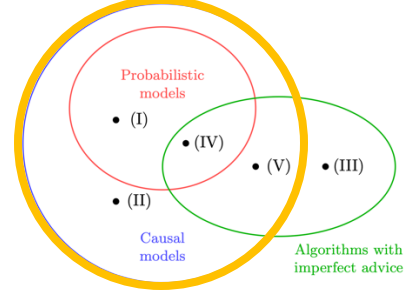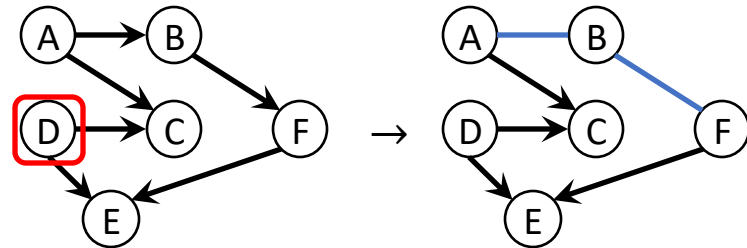
# Causal discovery via interventions



Ground truth
causal graph $\mathcal{G}^*$

Essential
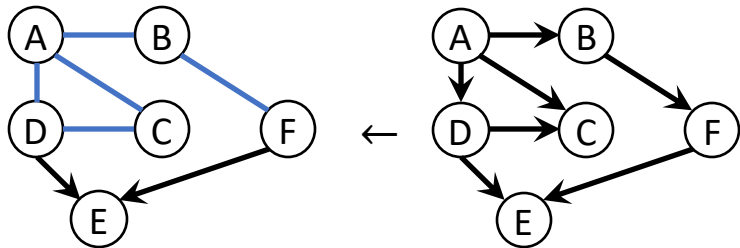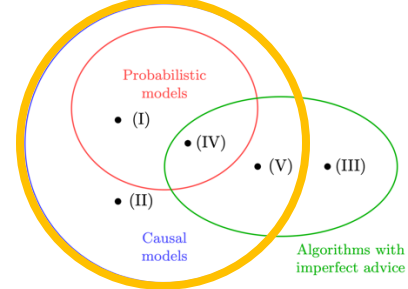graph $\mathcal{E}(\mathcal{G}^*)$

- Want: Recover $\mathcal{G}^*$ starting from partially oriented $\mathcal{E}(\mathcal{G}^*)$ from observational data

# Causal discovery via interventions



Ground truth causal graph $\mathcal{G}^*$

Essential graph $\mathcal{E}(\mathcal{G}^*)$

- Want: Recover $\mathcal{G}^*$ starting from partially oriented $\mathcal{E}(\mathcal{G}^*)$ from observational data
- Interventions reveal arc orientations (incident arcs + Meek rules)
- **Goal: Recover $\mathcal{G}^*$ using as few interventions as possible**

# Causal discovery via interventions



Ground truth
causal graph $\mathcal{G}^*$
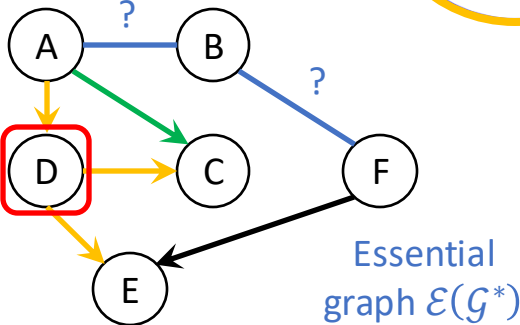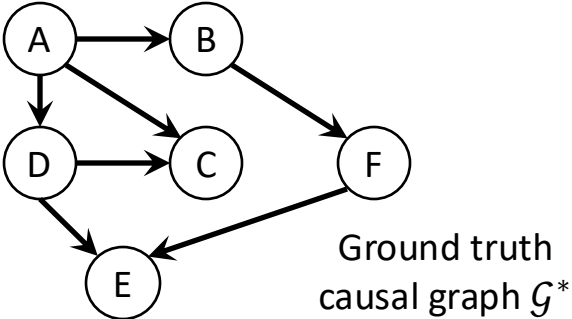
Essential
graph $\mathcal{E}(\mathcal{G}^*)$

- Want: Recover $\mathcal{G}^*$ starting from partially oriented $\mathcal{E}(\mathcal{G}^*)$ from observational data
- Interventions reveal arc orientations (incident arcs + Meek rules)
- **Goal: Recover $\mathcal{G}^*$ using as few interventions as possible**
- We have some results regarding how to design algorithms to perform optimal adaptive interventions under various scenarios [CSB22, CS23a, CGB23, CS23b, CS23c, CSU24]
- Insight: Reduce to graph / set cover problem with specialized causal operations

[CSB22] Davin Choo, Kirankumar Shiragur, Arnab Bhattacharyya. *Verification and search algorithms for causal DAGs*. Conference on Neural Information Processing Systems (NeurIPS), 2022.
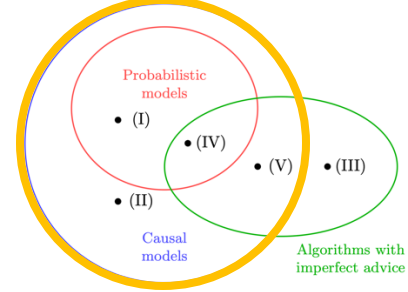[CS23a] Davin Choo, Kirankumar Shiragur. *Subset verification and search algorithms for causal DAGs*. International Conference on Artificial Intelligence and Statistics (AISTATS), 2023.
[CGB23] Davin Choo, Themistoklis Gouleakis, Arnab Bhattacharyya. *Active causal structure learning with advice*. International Conference on Machine Learning (ICML), 2023.
[CS23b] Davin Choo, Kirankumar Shiragur. *New metrics and search algorithms for weighted causal DAGs*. International Conference on Machine Learning (ICML), 2023.
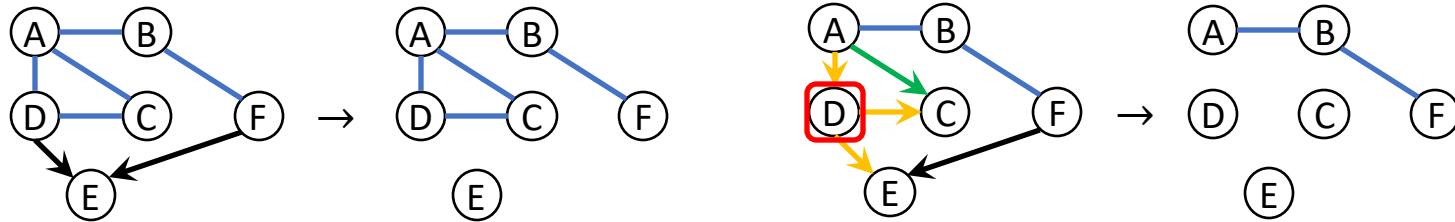[CS23c] Davin Choo, Kirankumar Shiragur. *Adaptivity Complexity for Causal Graph Discovery*. Conference on Uncertainty in Artificial Intelligence (UAI), 2023.
[CSU24] Davin Choo, Kirankumar Shiragur, Caroline Uhler. *Causal discovery under off-target interventions*. International Conference on Artificial Intelligence and Statistics (AISTATS), 2024.
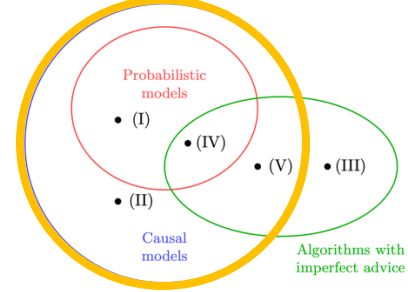
# A glimpse of [CSB22]



- **Insight: Frame as graph problem with causal operations**
- Known facts and observations (say n vertices)
  - Remove directed edges in essential graph → chordal graph $G$
  - If $G$ has no (undirected) edges, then whole graph is oriented
  - Intervention on vertex v → Orient all edges incident to v (possibly more)


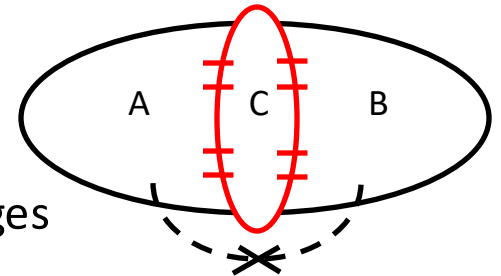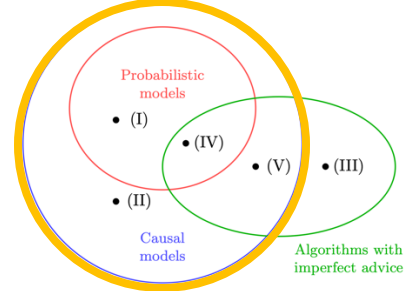
Essential graphs from earlier slides

# A glimpse of [CSB22]



- Insight: Frame as graph problem with causal operations
- Known facts and observations (say n vertices)
  - Remove directed edges in essential graph → chordal graph $G$
  - If $G$ has no (undirected) edges, then whole graph is oriented
  - Intervention on vertex v → Orient all edges incident to v (possibly more)
- Chordal graph separators [GRE84]
  - $|A|, |B| \leq \frac{|G|}{2}$ and $C$ is a clique, i.e., $|C| \leq \omega(G)$
  - Intervene on vertices in $C$ one by one
  - Repeat $O(\log n)$ times → $G$ will have no more edges



- We also show that this is optimal in worst case

[GRE83] John R. Gilbert, Donald J. Rose, Anders Edenbrandt. *A Separator Theorem for Chordal Graphs*. SIAM Journal on Algebraic Discrete Methods, 1984.
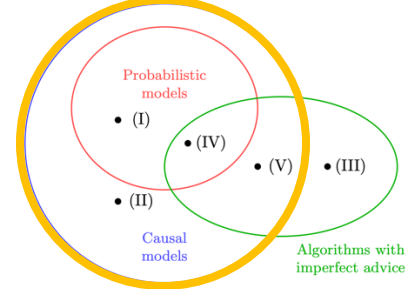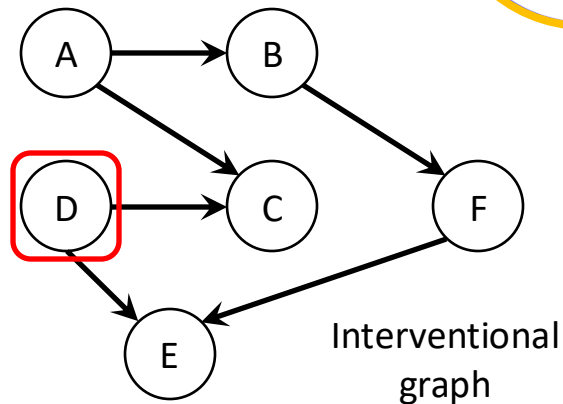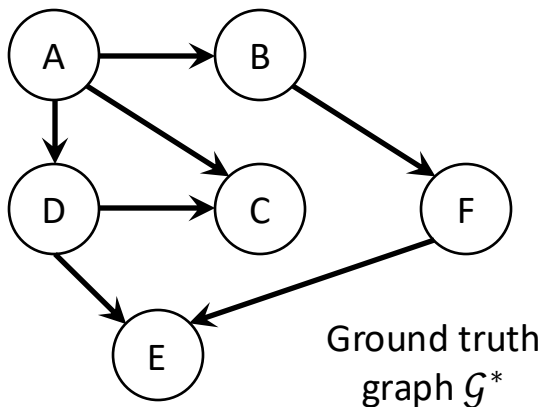
# (II): Causal models



- Two fundamental problems in causal inference
  - Causal graph discovery: Recover true causal graph $\mathcal{G}^*$
    - Even with infinite observational data, can only determine causal graph up to some equivalence class where all conditional independence relations agree
    - Make distributional/structural assumptions or perform interventions/experiments!
  - Causal effect estimation: Estimate $\mathcal{P}(Y = y \mid do(X = x))$
    - Typically, a 2-stage process: learn $\mathcal{G}^*$, then apply closed-form formulas

# Causal identification (the 2$^{nd}$ step)



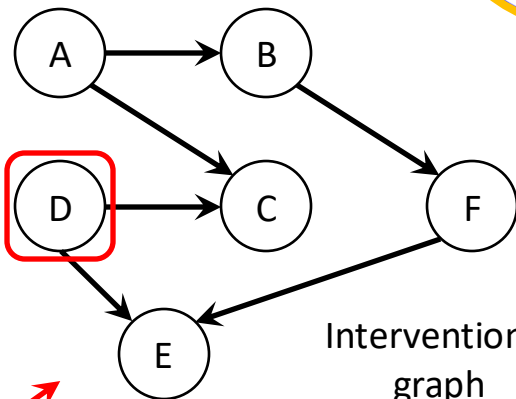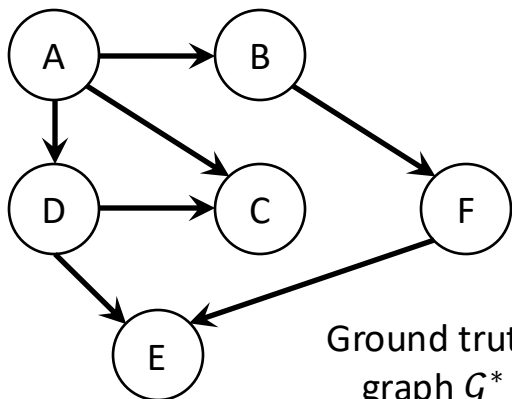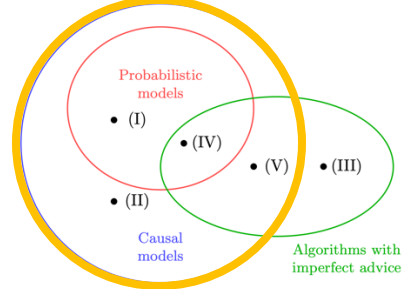Ground truth graph $\mathcal{G}^*$

Interventional graph

$$\mathcal{P}(E = e \mid do(D = d^*))$$

Interventional query

What is probability of $E = e$ when we fix $D = d^*$?

# Causal identification (the 2<sup>nd</sup> step)



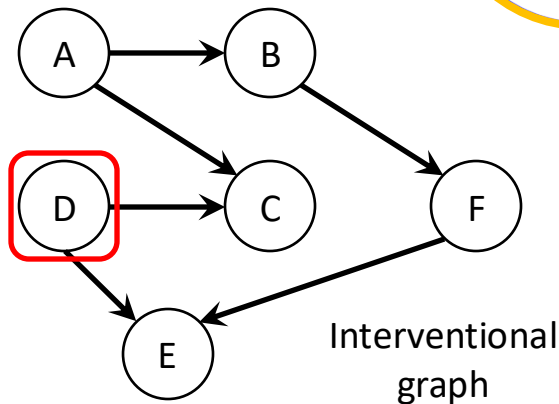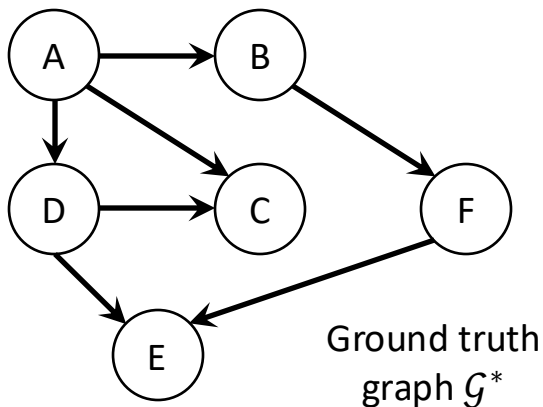Ground truth graph $\mathcal{G}^*$

Interventional graph

$$\mathcal{P}(E = e \mid do(D = d^*)) = \mathcal{P}(e \mid do(d^*)) \neq \mathcal{P}(e \mid d^*) \text{ in general}$$

Interventional query
What is probability of $E = e$ when we fix $D = d^*$?

Need to draw samples from interventional graph, i.e., perform experiment and measure

# Causal identification (the 2nd step)

Ground truth graph $\mathcal{G}^*$

Interventional graph
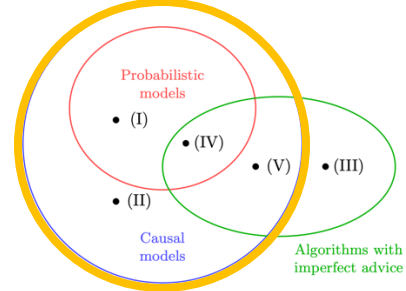
Because of structure of $\mathcal{G}^*$

$$\mathcal{P}(E = e \mid do(D = d^*)) = \mathcal{P}(e \mid do(d^*)) = \int \mathcal{P}(e \mid d^*, a) \cdot \mathcal{P}(a) \, da$$

Interventional query
What is probability of $E = e$ when we fix $D = d^*$?

Just observational terms!

# (II): Causal models



- ## Two fundamental problems in causal inference
  - Causal graph discovery: Recover true causal graph $\mathcal{G}^*$
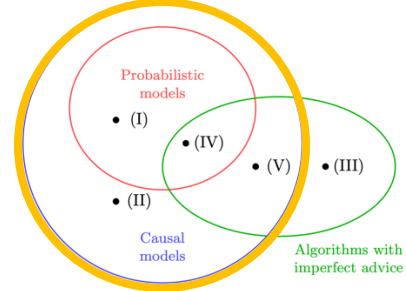    - Even with infinite observational data, can only determine causal graph up to some equivalence class where all conditional independence relations agree
    - Make distributional/structural assumptions or perform interventions/experiments!
  - Causal effect estimation: Estimate $\mathcal{P}(Y = y \,|\, do(X = x))$
    - Typically, a 2-stage process: learn $\mathcal{G}^*$, then apply closed-form formulas
    - [CSBS24] This is suboptimal as it may require strong assumptions and a lot of samples
    - Insight: "weak edges" shouldn't affect much for PAC-style results

[CSBS24] Davin Choo, Chandler Squires, Arnab Bhattacharyya, and David Sontag. *Probably approximately correct high-dimensional causal effect estimation given a valid adjustment set*. Under submission, 2024.

# A glimpse of [C̲SBS24]



Insight: "weak edges" shouldn't affect much for PAC-style results

- Suppose we draw observational samples from this causal graph of binary variables and wish to estimate interventional effect $\mathcal{P}(Y = y \mid do(X = x^*))$

- Let's estimate $\mathcal{P}(y \mid do(x^*))$ via $\sum_s \mathcal{P}(y|x^*, z)\mathcal{P}(z)$ for some subset $Z \subseteq V$

$V$ excludes $X$ and $Y$
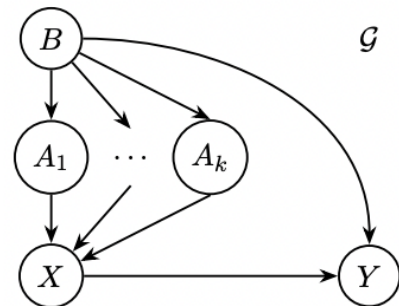
# A glimpse of [CSBS24]



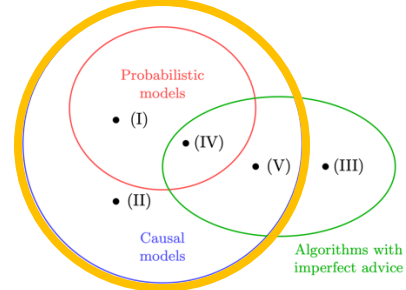Insight: "weak edges" shouldn't affect much for PAC-style results

- Suppose we draw observational samples from this causal graph of binary variables and wish to estimate interventional effect $\mathcal{P}(Y = y \mid do(X = x^*))$

- Let's estimate $\mathcal{P}(y \mid do(x^*))$ via $\sum_s \mathcal{P}(y|x^*, z)\mathcal{P}(z)$ for some subset $Z \subseteq V$
  - Valid when $Z$ is $\{B, A_1, \dots, A_k\}$ or $\{A_1, \dots, A_k\}$ or $\{B\}$, for <u>any</u> underlying $\mathcal{P}$
  - $\{B\}$ is the best: smaller set = less samples for an accurate estimate
  - But… we don't know the graph!

$V$ excludes $X$ and $Y$



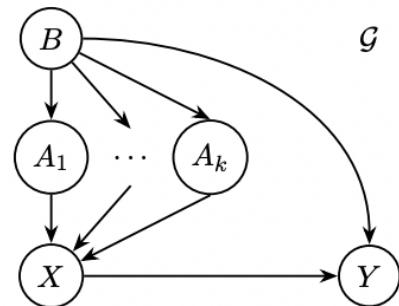17

# A glimpse of [C̲SBS24]

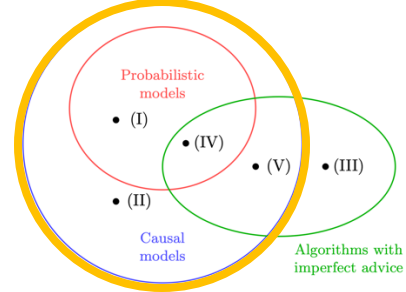

<span style="color:orange">Insight: "weak edges" shouldn't affect much for PAC-style results</span>

- Suppose we draw observational samples from this causal graph of binary variables and wish to estimate interventional effect $\mathcal{P}(Y = y \mid do(X = x^*))$

- Let's estimate $\mathcal{P}(y \mid do(x^*))$ via $\sum_s \mathcal{P}(y|x^*, z)\mathcal{P}(z)$ for some subset $Z \subseteq V$
  - Valid when $Z$ is $\{B, A_1, \ldots, A_k\}$ or $\{A_1, \ldots, A_k\}$ or $\{B\}$, for <u>any</u> underlying $\mathcal{P}$
  - $\{B\}$ is the best: smaller set = less samples for an accurate estimate
  - But… we don't know the graph!

- What CI tests + do-calculus that will validate the estimate?
  - "Markov blanket": $X \perp\!\!\!\perp S\backslash V \mid S \rightarrow$ Get $S = \{A_1, \ldots, A_k\}$
  - "Screening set": $Y \perp\!\!\!\perp S\backslash S' \mid X \cup S'$ and $X \perp\!\!\!\perp S'\backslash S \mid S \rightarrow$ Get $S' = \{B\}$

<span style="color:blue">$V$ excludes $X$ and $Y$</span>



17
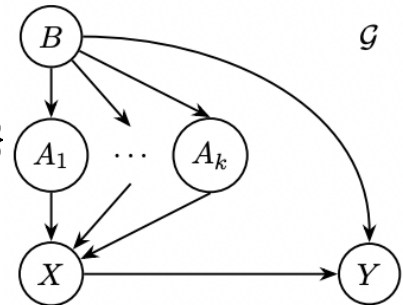
# A glimpse of [C̲SBS24]



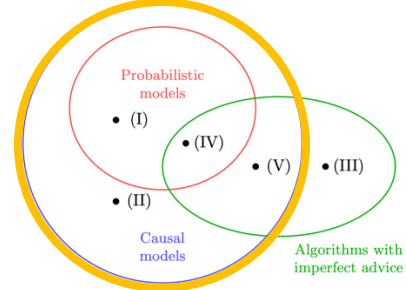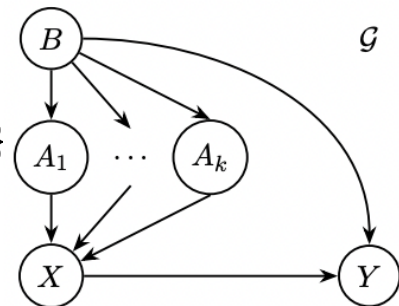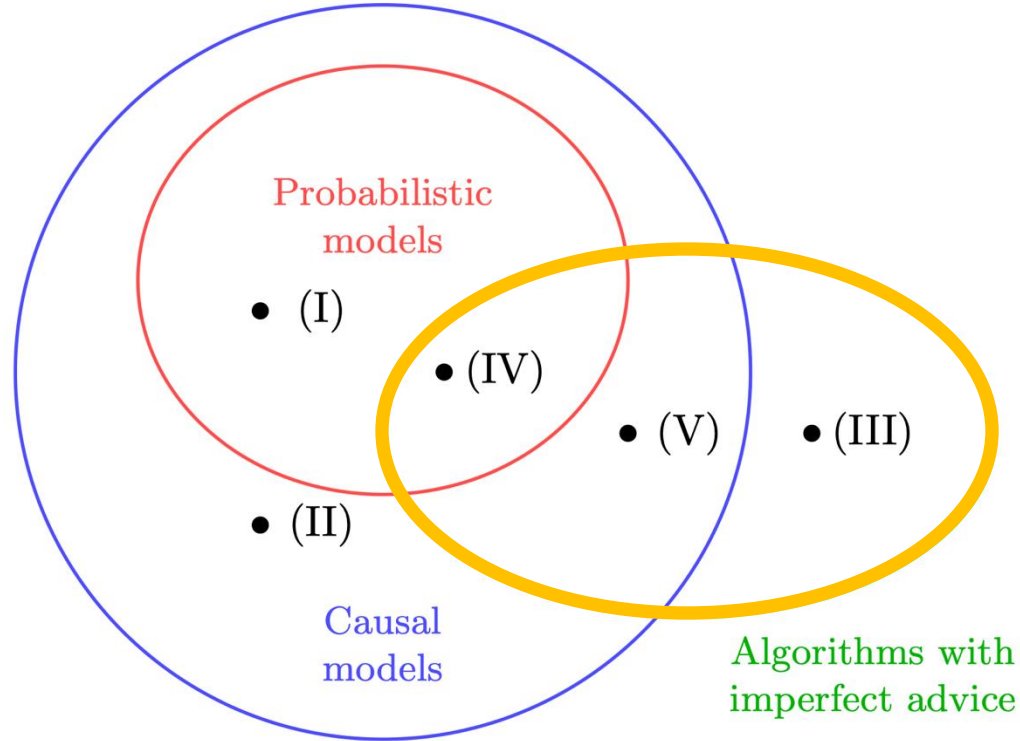Insight: "weak edges" shouldn't affect much for PAC-style results

- Suppose we draw observational samples from this causal graph of binary variables and wish to estimate interventional effect $\mathcal{P}(Y = y \mid do(X = x^*))$

- Let's estimate $\mathcal{P}(y \mid do(x^*))$ via $\sum_s \mathcal{P}(y|x^*, z)\mathcal{P}(z)$ for some subset $Z \subseteq V$
  - Valid when $Z$ is $\{B, A_1, \dots, A_k\}$ or $\{A_1, \dots, A_k\}$ or $\{B\}$, for <u>any</u> underlying $\mathcal{P}$
  - $\{B\}$ is the best: smaller set = less samples for an accurate estimate
  - But… we don't know the graph!

- What CI tests + do-calculus that will validate the estimate?
  - "Markov blanket": $X \perp\!\!\!\perp S\backslash \text{V} \mid S \rightarrow$ Get $S = \{A_1, \dots, A_k\}$
  - "Screening set": $Y \perp\!\!\!\perp S\backslash S' \mid X \cup S'$ and $X \perp\!\!\!\perp S'\backslash S \mid S \rightarrow$ Get $S' = \{B\}$
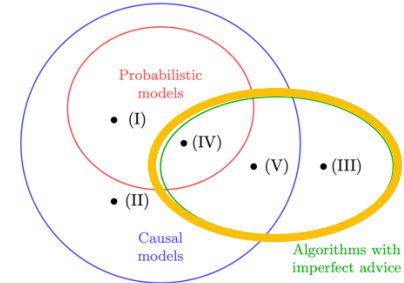  - Approximate conditional independence test $\rightarrow$ PAC estimate

$V$ excludes $X$ and $Y$
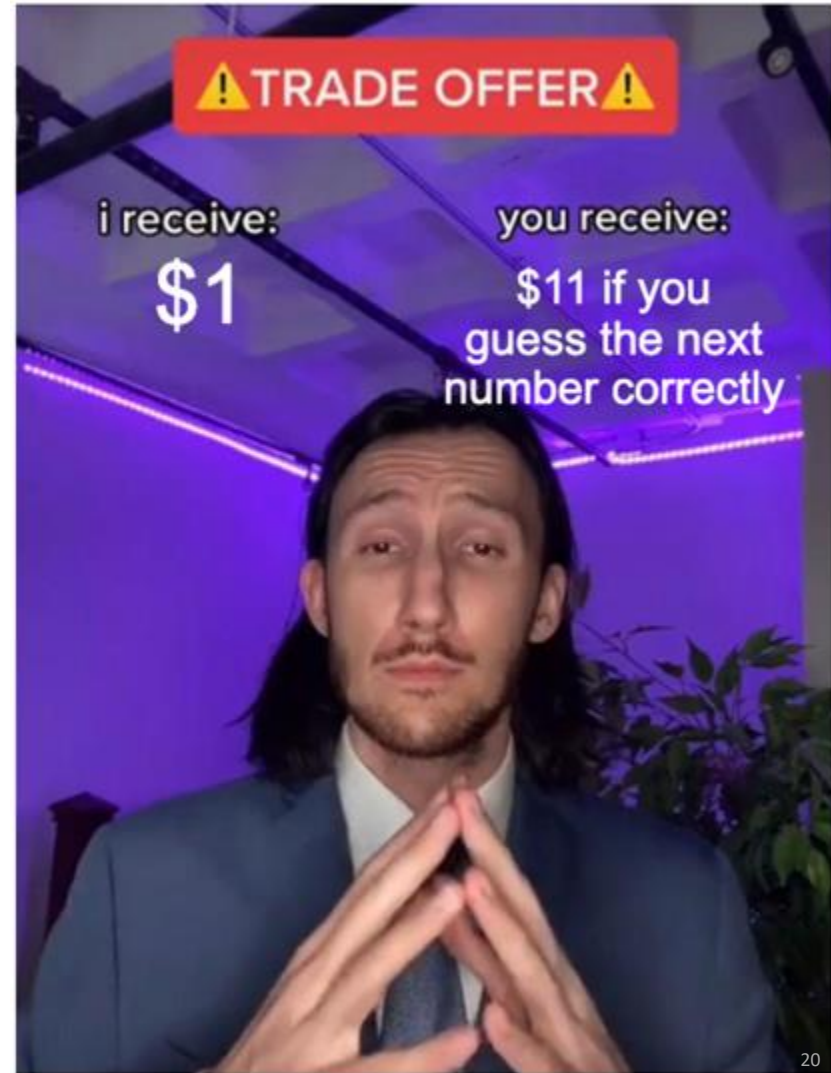


17

# Main themes explored in my PhD thesis

# (III/IV/V): Algorithms with advice



- Two key performance measures
  - Consistency: If advice is "perfect", how good are things?
  - Robustness: If advice is "garbage", how bad are things?
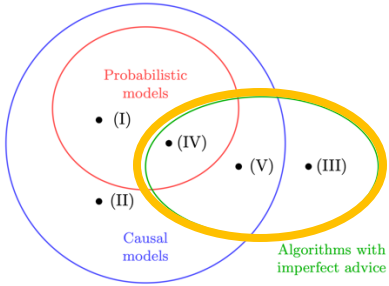- Challenge: We don't know how good the given advice is a priori!

# Detour: Let's make a deal

- There are 10 numbers in the universe U = {0, 1, 2, 3, 4, 5, 6, 7, 8, 9}

- There is an underlying process $\mathcal{P}$ that generates i.i.d. samples from U
  - i.e., We can observe a sequence such as 1, 6, 3, 6, 2, 8, 0, 3, 9, 5, 4, …

- What property of $\mathcal{P}$ will make this deal <u>profitable in expectation</u>?



⚠️TRADE OFFER⚠️

i receive:

$1

you receive:

$11 if you guess the next number correctly

# Detour: Property testing land

- How to test if $\mathcal{P}$ is the uniform distribution over U?
  - Say, we only care about constant success probability (can be amplified)

- Learning a $\varepsilon$-close $\widehat{\mathcal{P}}$ then check: $\Theta\left(\frac{|U|}{\varepsilon^2}\right)$ i.i.d. samples from $\mathcal{P}$

- Uniformity testing requires $\Theta\left(\frac{\sqrt{|U|}}{\varepsilon^2}\right)$ i.i.d. samples from $\mathcal{P}$
  - If $\mathcal{P}$ is uniform, output YES w.p. $\geq \frac{2}{3}$
  - If $\mathcal{P}$ is $\varepsilon$-far from uniform, output NO w.p. $\geq \frac{2}{3}$
  - Many existing proofs for this bound. E.g., look at collisions in samples

  $\longrightarrow$ Allowed to output arbitrarily if not uniform, yet not "far from uniform"

- See also [Can22] for an excellent property testing survey

[Can22] Clément L. Canonne. *Topics and Techniques in Distribution Testing: A Biased but Representative Sample*. Foundations and Trends® in Communications and Information Theory, 2022.
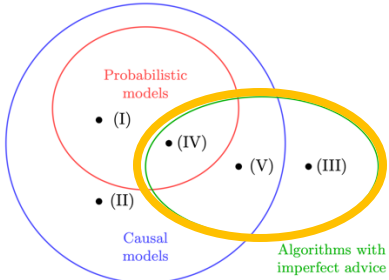
# (III/IV/V): Algorithms with advice



- Two key performance measures
  - Consistency: If advice is "perfect", how good are things?
  - Robustness: If advice is "garbage", how bad are things?

- Challenge: We don't know how good the given advice is a priori!

- Insight: "Testing can be cheaper than learning" → TestAndAct
  - [CGLB24] TestAndMatch: Improve competitive ratio of online bipartite matching (III)
  - [BCGG24] TestAndOptimize: Improve sample complexity of learning multivariate Gaussians (IV)
  - [CGB23] TestAndSubsetSearch: Reduce num of interventions required for causal graph discovery (V)

[CGLB24] Davin Choo, Themistoklis Gouleakis, Chun Kai Ling, and Arnab Bhattacharyya. *Online bipartite matching with imperfect advice*. International Conference on Machine Learning (ICML), 2024.
[BCGG24] Arnab Bhattacharyya, Davin Choo, Philips George John, and Themistoklis Gouleakis. *Learning multivariate Gaussians with imperfect advice*. Under submission, 2024.
[CGB23] Davin Choo, Themistoklis Gouleakis, Arnab Bhattacharyya. *Active causal structure learning with advice*. International Conference on Machine Learning (ICML), 2023.
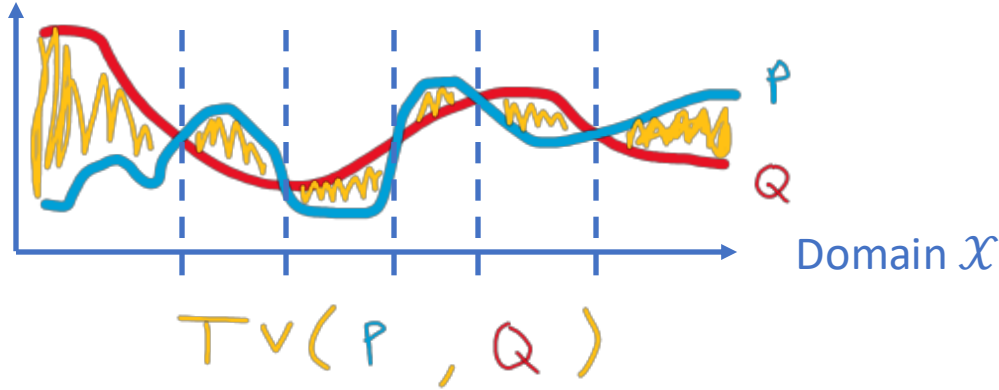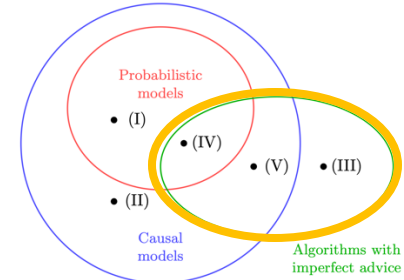
# A glimpse of [B**C**GG24]



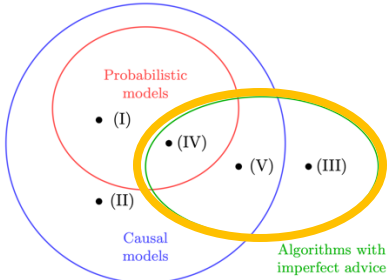Insight: "Testing can be cheaper than learning"

- Gaussian estimation with i.i.d. samples
  - Given sample access to some underlying distribution $\mathcal{P}$, produce $\widehat{\mathcal{P}}$ such that $\text{TV}(\mathcal{P}, \widehat{\mathcal{P}}) \leq \varepsilon$ with probability $\geq 1 - \delta$



Probability mass, i.e. area under curve sums to 1

Domain $\mathcal{X}$

$TV(P, Q)$

# A glimpse of [BCGG24]



Insight: "Testing can be cheaper than learning"

- Gaussian estimation with i.i.d. samples
  - Given sample access to some underlying distribution $\mathcal{P}$, produce $\hat{\mathcal{P}}$ such that $\text{TV}(\mathcal{P}, \hat{\mathcal{P}}) \leq \varepsilon$ with probability $\geq 1 - \delta$

- Useful to invoke Pinsker's inequality: $\text{TV}(\mathcal{P}, \hat{\mathcal{P}}) \leq \sqrt{\frac{1}{2} \cdot \text{KL}(\mathcal{P}, \hat{\mathcal{P}})}$

- For multivariate Gaussians over $\mathbb{R}^d$,

$$\text{KL}\big(N(\boldsymbol{\mu}_{\mathcal{P}}, \boldsymbol{\Sigma}_{\mathcal{P}}), N(\boldsymbol{\mu}_{\mathcal{Q}}, \boldsymbol{\Sigma}_{\mathcal{Q}})\big) = \frac{1}{2} \cdot \left[\text{Tr}(\boldsymbol{\Sigma}_{\mathcal{Q}}^{-1}\boldsymbol{\Sigma}_{\mathcal{P}}) - d + \ln\left(\frac{\det \boldsymbol{\Sigma}_{\mathcal{Q}}}{\det \boldsymbol{\Sigma}_{\mathcal{P}}}\right)\right]$$
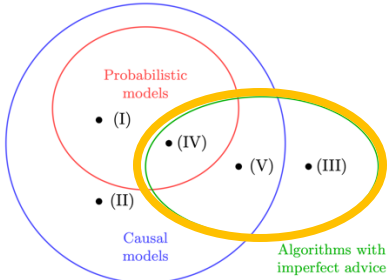
- So, we just need to upper bound KL by $\varepsilon^2$
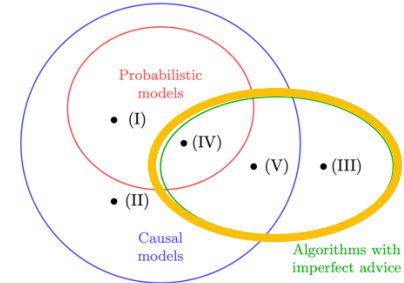
# A glimpse of [B**C**GG24]



Insight: "Testing can be cheaper than learning"

- Let's consider the simple identity covariance setting

$$\text{KL}\big(N(\boldsymbol{\mu}, \mathbf{I}_d), N(\widehat{\boldsymbol{\mu}}, \mathbf{I}_d)\big) = \frac{1}{2} \cdot \|\boldsymbol{\mu} - \widehat{\boldsymbol{\mu}}\|_2^2$$

Linear in dimension d

- Empirical estimator is optimal: need $\widetilde{\Theta}\left(\frac{d}{\varepsilon^2}\right)$ samples to get $\text{KL} \leq \varepsilon^2$

# A glimpse of [BCGG24]



Insight: "Testing can be cheaper than learning"

- Let's consider the simple identity covariance setting

$$\text{KL}\big(N(\boldsymbol{\mu}, \mathbf{I}_d), N(\widehat{\boldsymbol{\mu}}, \mathbf{I}_d)\big) = \frac{1}{2} \cdot \|\boldsymbol{\mu} - \widehat{\boldsymbol{\mu}}\|_2^2$$

Linear in dimension d

- Empirical estimator is optimal: need $\widetilde{\Theta}\left(\dfrac{d}{\varepsilon^2}\right)$ samples to get KL $\leq \varepsilon^2$

- Can we do better if someone proposes $\widetilde{\boldsymbol{\mu}}$ as advice?
  - If $\widetilde{\boldsymbol{\mu}} = \boldsymbol{\mu}$, then 0 samples needed, but we cannot blindly trust it
  - W.L.O.G., can treat $\widetilde{\boldsymbol{\mu}} = \mathbf{0}_d$ by pre-processing the samples accordingly
    - Given samples $\boldsymbol{y_1}, \ldots, \boldsymbol{y_n} \sim \mathcal{P}$, consider $(\boldsymbol{y_1} - \widetilde{\boldsymbol{\mu}}), \ldots, (\boldsymbol{y_n} - \widetilde{\boldsymbol{\mu}})$ instead
    - Once we obtain estimate $\widehat{\boldsymbol{\mu}}$, output $\widehat{\boldsymbol{\mu}} + \widetilde{\boldsymbol{\mu}}$ instead
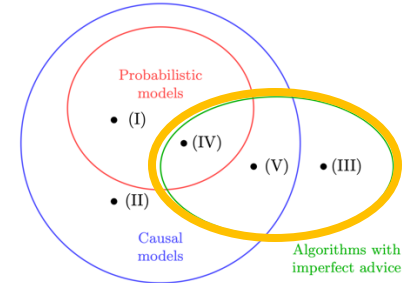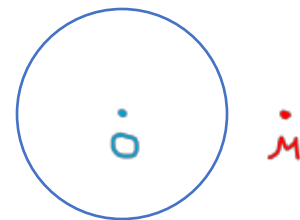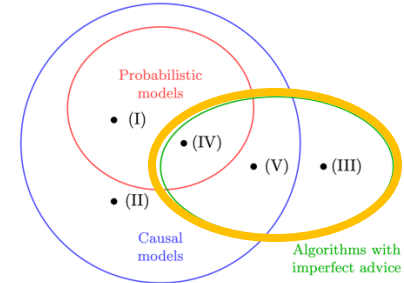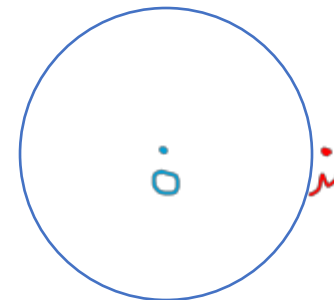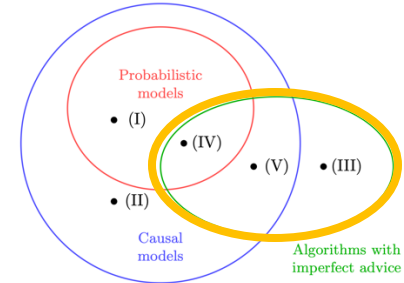
# A glimpse of [B<u>C</u>GG24]

Insight: "Testing can be cheaper than learning"

- High-level idea
  - Use sublinear tolerant testing + exponential search to find $r > 0$ s.t. $\frac{r}{2} \leq \|\boldsymbol{\mu}\|_2 \leq r$
  - Then, search within this radius to find a "good enough" $\widehat{\boldsymbol{\mu}}$
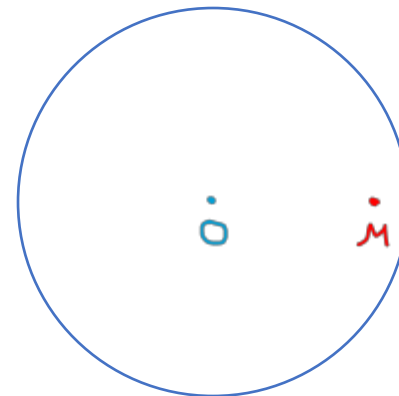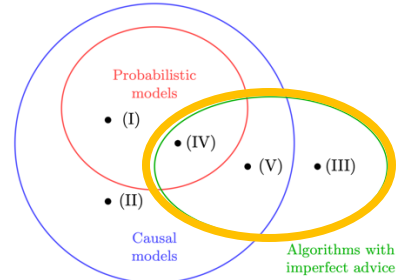
# A glimpse of [B**C**GG24]



Insight: "Testing can be cheaper than learning"

- High-level idea
    - Use sublinear tolerant testing + exponential search to find $r > 0$ s.t. $\frac{r}{2} \leq \|\boldsymbol{\mu}\|_2 \leq r$
    - Then, search within this radius to find a "good enough" $\widehat{\boldsymbol{\mu}}$

# A glimpse of [B<u>C</u>GG24]



Insight: "Testing can be cheaper than learning"

- High-level idea
    - Use sublinear tolerant testing + exponential search to find $r > 0$ s.t. $\frac{r}{2} \leq \|\boldsymbol{\mu}\|_2 \leq r$
    - Then, search within this radius to find a "good enough" $\widehat{\boldsymbol{\mu}}$

# A glimpse of [BCGG24]



Insight: "Testing can be cheaper than learning"

- High-level idea
  - Use sublinear tolerant testing + exponential search to find $r > 0$ s.t. $\frac{r}{2} \leq \|\boldsymbol{\mu}\|_2 \leq r$
  - Then, search within this radius to find a "good enough" $\widehat{\boldsymbol{\mu}}$

# A glimpse of [B**C**GG24]



Insight: "Testing can be cheaper than learning"

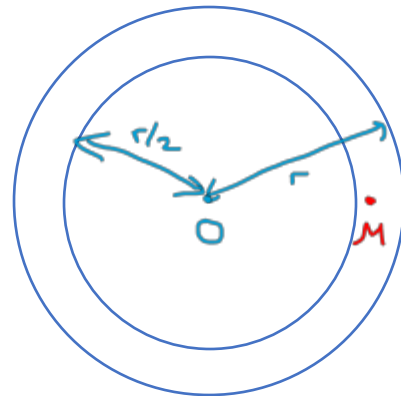Each test uses $\tilde{\mathcal{O}}\left(\frac{\sqrt{d}}{\varepsilon^2}\right)$, as compared to $\tilde{\Theta}\left(\frac{d}{\varepsilon^2}\right)$ for empirical estimator
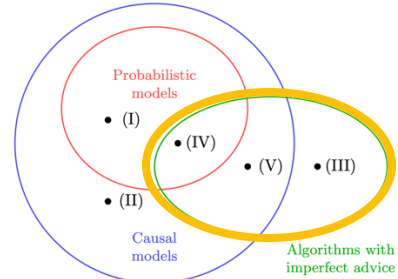
- High-level idea
  - Use sublinear tolerant testing + exponential search to find $r > 0$ s.t. $\frac{r}{2} \leq \|\boldsymbol{\mu}\|_2 \leq r$
  - Then, search within this radius to find a "good enough" $\hat{\boldsymbol{\mu}}$

# A glimpse of [B**C**GG24]

Each test uses $\tilde{\mathcal{O}}\left(\frac{\sqrt{d}}{\varepsilon^2}\right)$, as compared to $\tilde{\Theta}\left(\frac{d}{\varepsilon^2}\right)$ for empirical estimator

Insight: "Testing can be cheaper than learning"

- High-level idea
  - Use sublinear tolerant testing + exponential search to find $r > 0$ s.t. $\frac{r}{2} \leq \|\boldsymbol{\mu}\|_2 \leq r$
  - Then, search within this radius to find a "good enough" $\hat{\boldsymbol{\mu}}$
  - For technical reasons, we need to estimate $\|\boldsymbol{\mu}\|_1$ with some $\lambda$ instead
  - Then, using i.i.d. samples $\boldsymbol{y}_1, \dots, \boldsymbol{y}_n$ from $\mathcal{P}$, solve LASSO in poly time:

$$\hat{\boldsymbol{\mu}} = \operatorname*{argmin}_{\|\boldsymbol{\beta}\|_1 \leq r} \frac{1}{n} \sum_{i=1}^{n} \|\boldsymbol{y}_i - \boldsymbol{\beta}\|_2^2$$

  - When $\|\boldsymbol{\mu}\|_1$ is sufficiently small, our method provably uses $\tilde{o}\left(\frac{d}{\varepsilon^2}\right)$
  - Recall: Empirical estimator is optimal: need $\tilde{\Theta}\left(\frac{d}{\varepsilon^2}\right)$ samples to get KL $\leq \varepsilon^2$

# A glimpse of [B<u>C</u>GG24]



Each test uses $\tilde{\mathcal{O}}\left(\frac{\sqrt{d}}{\varepsilon^2}\right)$, as compared to
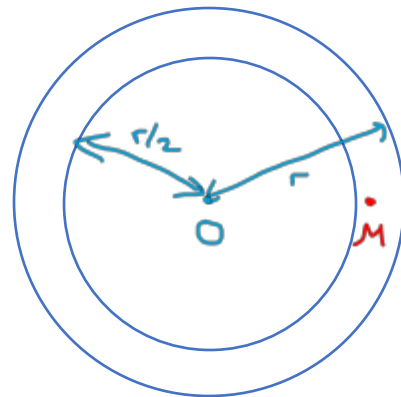
$\tilde{\Theta}\left(\frac{d}{\varepsilon^2}\right)$ for empirical estimator
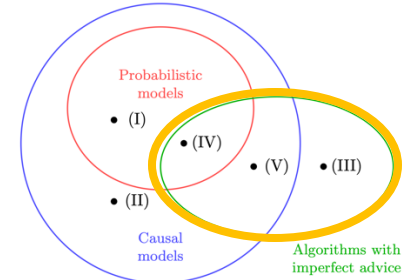
Insight: "Testing can be cheaper than learning"

- High-level idea
  - Use sublinear tolerant testing + exponential search to find $r > 0$ s.t. $\frac{r}{2} \le \|\boldsymbol{\mu}\|_2 \le r$
  - Then, search within this radius to find a "good enough" $\hat{\boldsymbol{\mu}}$
  - For technical reasons, we need to estimate $\|\boldsymbol{\mu}\|_1$ with some $\lambda$ instead
  - Then, using i.i.d. samples $\boldsymbol{y}_1, \dots, \boldsymbol{y}_n$ from $\mathcal{P}$, solve LASSO in poly time:

$$\hat{\boldsymbol{\mu}} = \operatorname{argmin}_{\|\boldsymbol{\beta}\|_1 \le r} \frac{1}{n}\sum_{i=1}^{n}\|\boldsymbol{y}_i - \boldsymbol{\beta}\|_2^2$$

  - When $\|\boldsymbol{\mu}\|_1$ is sufficiently small, our method provably uses $\tilde{o}\left(\frac{d}{\varepsilon^2}\right)$
  - Recall: Empirical estimator is optimal: need $\tilde{\Theta}\left(\frac{d}{\varepsilon^2}\right)$ samples to get KL $\le \varepsilon^2$

- Remarks
  - Small $\|\boldsymbol{\mu}\|_1$ here actually means small $\|\boldsymbol{\mu} - \tilde{\boldsymbol{\mu}}\|_1$ due to the pre-processing WLOG
  - We also need additional modifications tricks such as partitioning $\boldsymbol{\mu}$ into different coordinates to estimate $\|\boldsymbol{\mu}\|_1$, etc.
  - Similar idea work when the multivariate Gaussian has non-identity covariance matrix, but we use SDP instead of LASSO
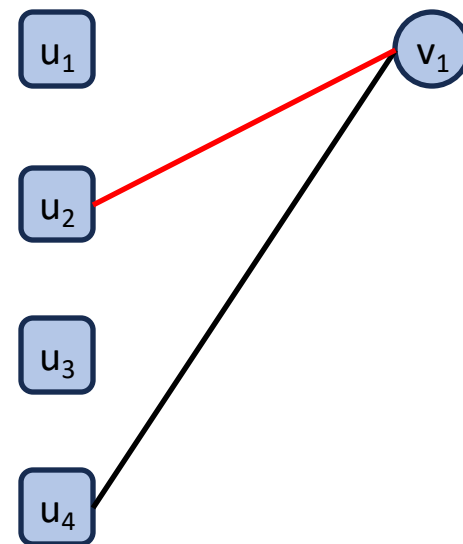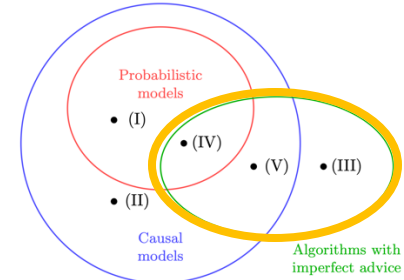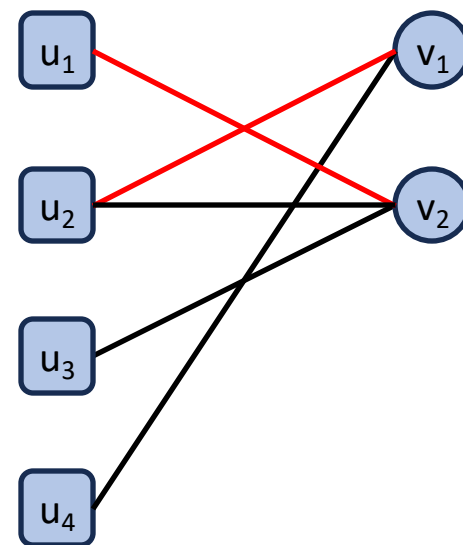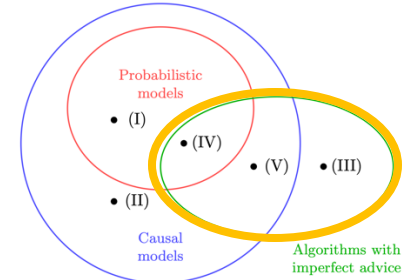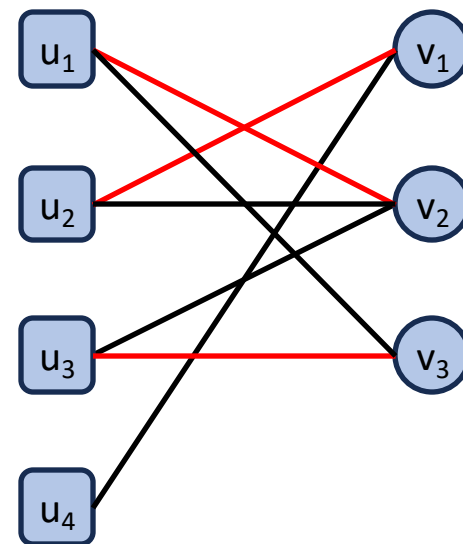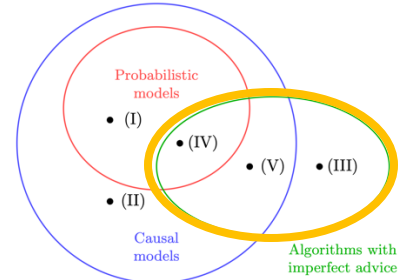
# A glimpse of [CGLB24]

Insight: "Testing can be cheaper than learning"

- Online bipartite matching
    - Offline set $U = \{u_1, \ldots, u_n\}$ fixed and known
    - Online set $V = \{v_1, \ldots, v_n\}$ arrive one by one
    - When an online vertex $v_i$ arrives
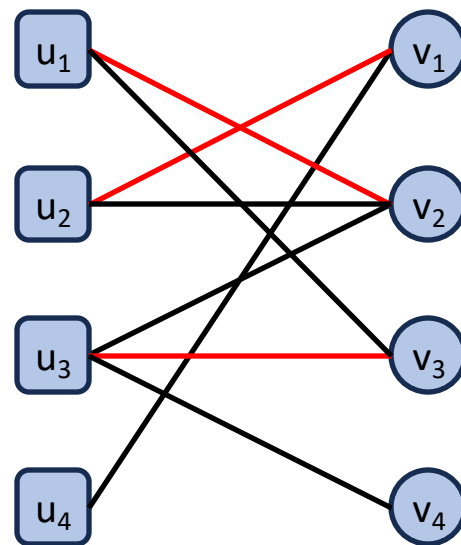        - $N(v_i)$ are revealed and we make <u>irrevocable</u> decision
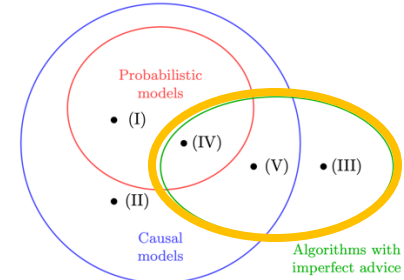
# A glimpse of [CGLB24]



Insight: "Testing can be cheaper than learning"

- Online bipartite matching
  - Offline set $U = \{u_1, \dots, u_n\}$ fixed and known
  - Online set $V = \{v_1, \dots, v_n\}$ arrive one by one
  - When an online vertex $v_i$ arrives
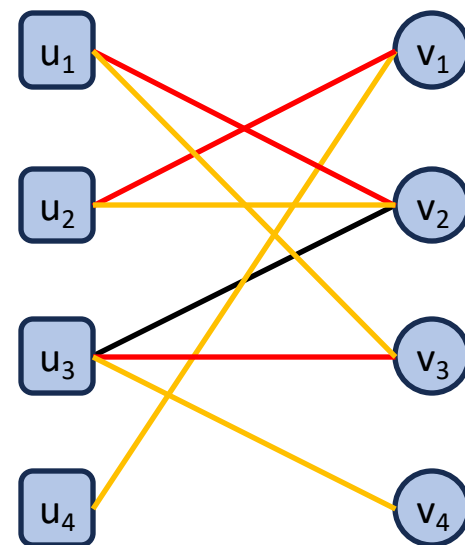    - $N(v_i)$ are revealed and we make <u>irrevocable</u> decision

# A glimpse of [CGLB24]



Insight: "Testing can be cheaper than learning"

- Online bipartite matching
  - Offline set $U = \{u_1, \dots, u_n\}$ fixed and known
  - Online set $V = \{v_1, \dots, v_n\}$ arrive one by one
  - When an online vertex $v_i$ arrives
    - $N(v_i)$ are revealed and we make <u>irrevocable</u> decision
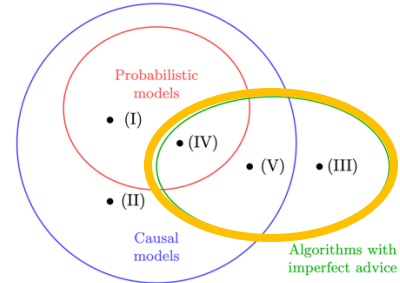
# A glimpse of [CGLB24]



Insight: "Testing can be cheaper than learning"

- Online bipartite matching
  - Offline set $U = \{u_1, \ldots, u_n\}$ fixed and known
  - Online set $V = \{v_1, \ldots, v_n\}$ arrive one by one
  - When an online vertex $v_i$ arrives
    - $N(v_i)$ are revealed and we make <u>irrevocable</u> decision

# A glimpse of [CGLB24]



Insight: "Testing can be cheaper than learning"

- Online bipartite matching
  - Offline set $U = \{u_1, \dots, u_n\}$ fixed and known
  - Online set $V = \{v_1, \dots, v_n\}$ arrive one by one
  - When an online vertex $v_i$ arrives
    - $N(v_i)$ are revealed and we make <u>irrevocable</u> decision
  - Final offline graph $G^* = (U \cup V, E)$
    - $E = N(v_1) \cup \dots \cup N(v_n)$
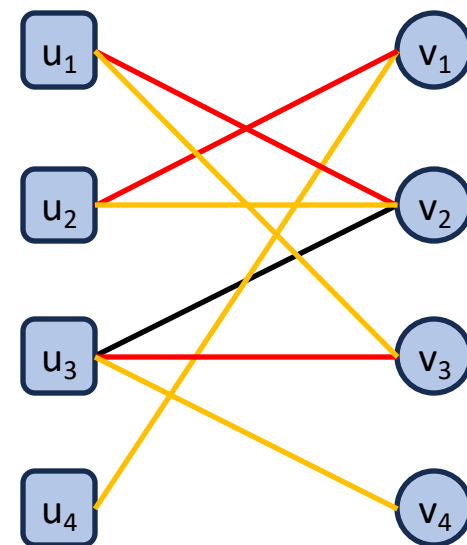    - Maximum matching $M^* \subseteq E$ of size $|M^*| = n^* \leq n$
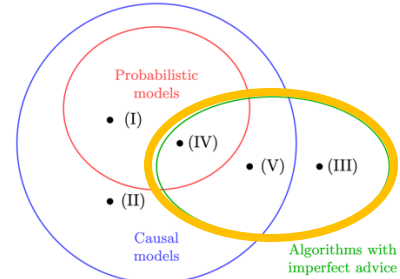
# A glimpse of [CGLB24]



Insight: "Testing can be cheaper than learning"

- Online bipartite matching
  - Offline set $U = \{u_1, \dots, u_n\}$ fixed and known
  - Online set $V = \{v_1, \dots, v_n\}$ arrive one by one
  - When an online vertex $v_i$ arrives
    - $N(v_i)$ are revealed and we make <u>irrevocable</u> decision
  - Final offline graph $G^* = (U \cup V, E)$
    - $E = N(v_1) \cup \cdots \cup N(v_n)$
    - Maximum matching $M^* \subseteq E$ of size $|M^*| = n^* \leq n$
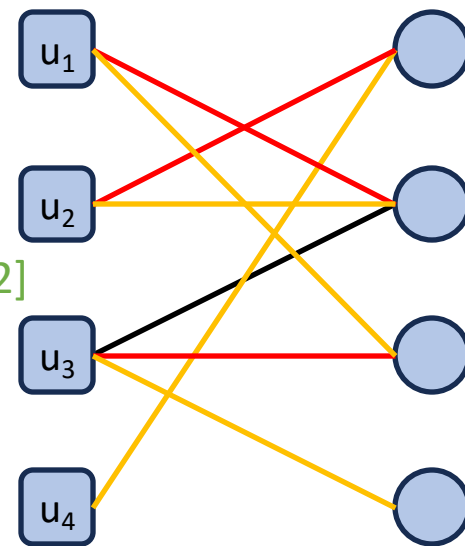  - Goal: Produce $M$ maximizing competitive ratio $\dfrac{|M|}{|M^*|}$



Here, the ratio is 3/4

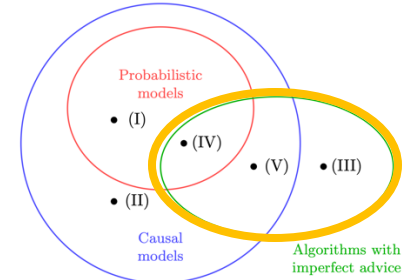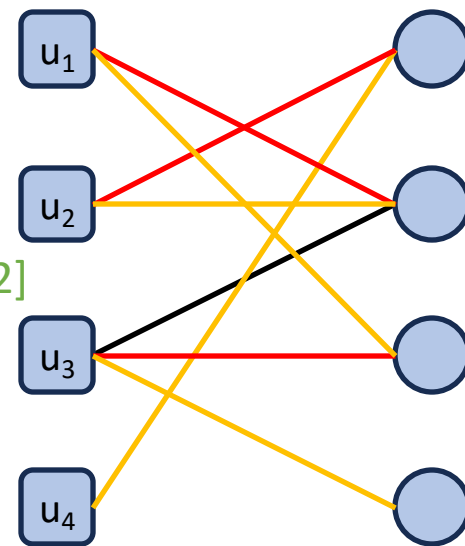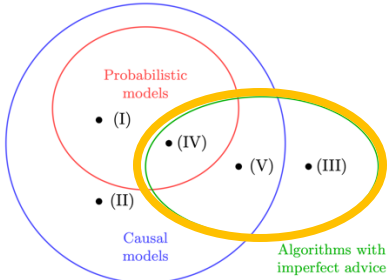# A glimpse of [CGLB24]



Insight: "Testing can be cheaper than learning"

- Online bipartite matching with <u>random arrival</u>
  - Still worst-case final graph $G^*$
  - Online vertex sequence is random permutation of V
  - Ranking achieves comp. ratio of 0.696 [MY11]
  - <u>No</u> algorithm cannot beat comp. ratio of 0.823 [MGS12]

- What we show
  - Advice = Prediction $\tilde{G}$ of $G^*$
  - When advice perfect ($\tilde{G} = G^*$), get comp. ratio 1
  - When advice bad, we get $\approx \beta$ ($0.696 \leq \beta \leq 0.823$)
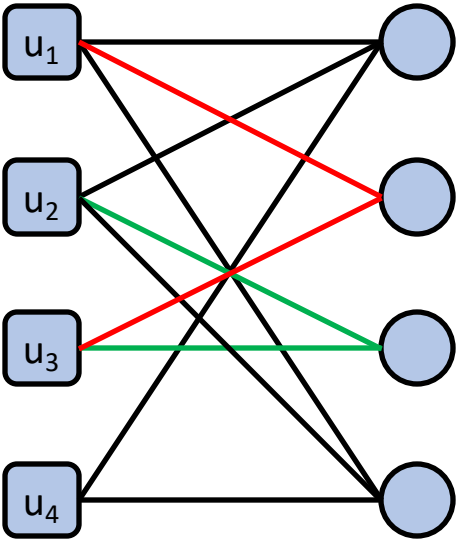
[MY11] Mohammad Mahdian and Qiqi Yan. *Online Bipartite Matching with Random Arrivals: An Approach Based on Strongly Factor-Revealing LPs*. Symposium on Theory of Computing (STOC), 2011
[MGS12] Vahideh H Manshadi, Shayan Oveis Gharan, and Amin Saberi. *Online stochastic matching: Online actions based on offline statistics*. Mathematics of Operations Research, 2012

# A glimpse of [CGLB24]



**Insight: "Testing can be cheaper than learning"**

- Online bipartite matching with <u>random arrival</u>
  - Still worst-case final graph $G^*$
  - Online vertex sequence is random permutation of V
  - Ranking achieves comp. ratio of 0.696 [MY11]
  - <u>No</u> algorithm cannot beat comp. ratio of 0.823 [MGS12]

- What we show
  - Advice = Prediction $\tilde{G}$ of $G^*$
  - When advice perfect ($\tilde{G} = G^*$), get comp. ratio 1
  - When advice bad, we get $\approx \beta$ ($0.696 \leq \beta \leq 0.823$)

Say, Baseline achieves this

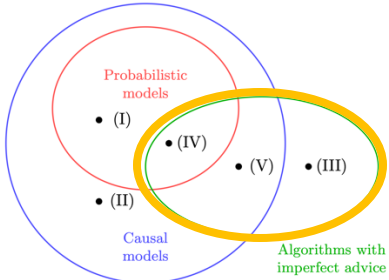[MY11] Mohammad Mahdian and Qiqi Yan. *Online Bipartite Matching with Random Arrivals: An Approach Based on Strongly Factor-Revealing LPs*. Symposium on Theory of Computing (STOC), 2011
[MGS12] Vahideh H Manshadi, Shayan Oveis Gharan, and Amin Saberi. *Online stochastic matching: Online actions based on offline statistics*. Mathematics of Operations Research, 2012

# A glimpse of [CGLB24]



Insight: "Testing can be cheaper than learning"

- Realized type counts as advice



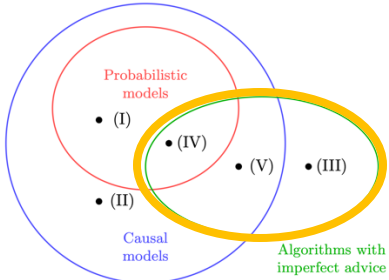| Type | $c^*$ |
|---|---|
| $\{u_1, u_2, u_4\}$ | 2 |
| $\{u_1, u_3\}$ | 1 |
| $\{u_2, u_3\}$ | 1 |
| $2^U \setminus T^*$ | 0 |

$T^*$

# A glimpse of [C̲GLB24]



Insight: "Testing can be cheaper than learning"

- Mimic algorithm: Fix arbitrary maximum matching $\widehat{M}$ defined by $\hat{c}$ and try to follow it as much as possible. If unable, leave unmatched
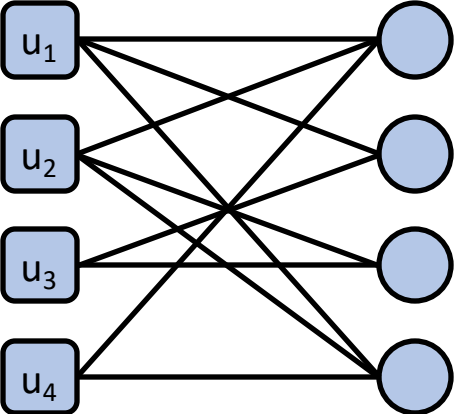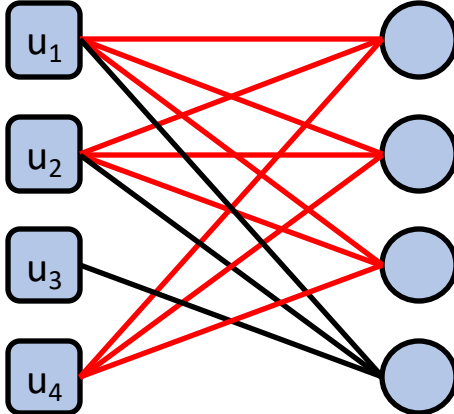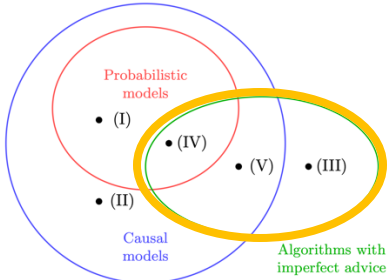


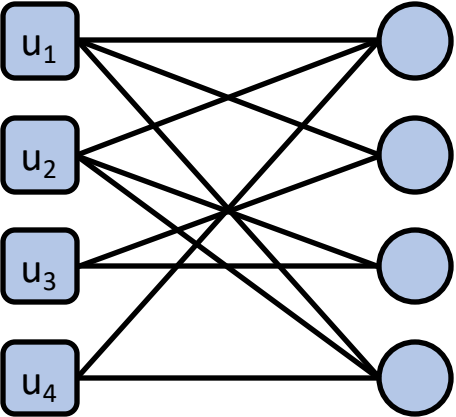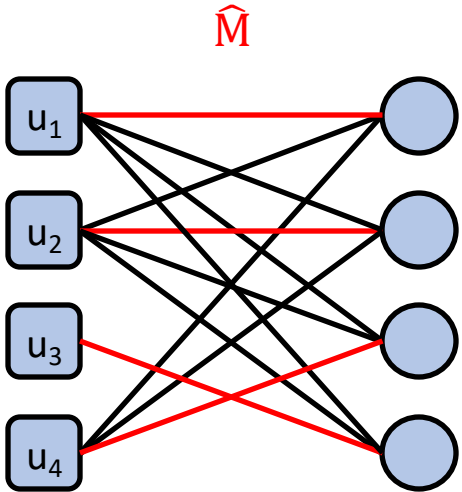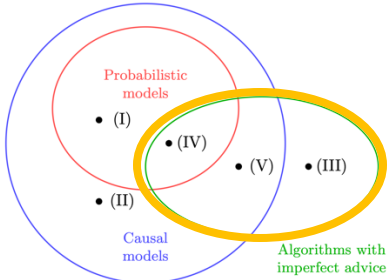| Type | $c^*$ | $\hat{c}$ |
|---|---|---|
| $\{u_1, u_2, u_4\}$ | 2 | 3 |
| $\{u_1, u_3\}$ | 1 | 0 |
| $\{u_2, u_3\}$ | 1 | 0 |
| $\{u_1, u_2, u_3\}$ | 0 | 1 |

# A glimpse of [C̲GLB24]



Insight: "Testing can be cheaper than learning"

- Mimic algorithm: Fix arbitrary maximum matching $\widehat{M}$ defined by $\hat{c}$ and try to follow it as much as possible. If unable, leave unmatched



| Type | $c^*$ | $\hat{c}$ |
|---|---|---|
| $\{u_1, u_2, u_4\}$ | 2 | 3 |
| $\{u_1, u_3\}$ | 1 | 0 |
| $\{u_2, u_3\}$ | 1 | 0 |
| $\{u_1, u_2, u_3\}$ | 0 | 1 |

# A glimpse of [CGLB24]



Insight: "Testing can be cheaper than learning"

- Mimic algorithm: Fix arbitrary maximum matching $\widehat{M}$ defined by $\hat{c}$ and try to follow it as much as possible. If unable, leave unmatched
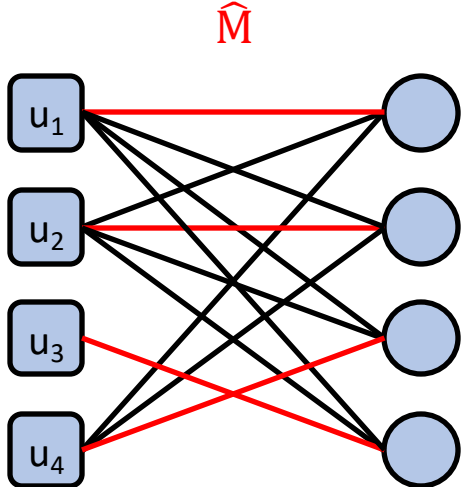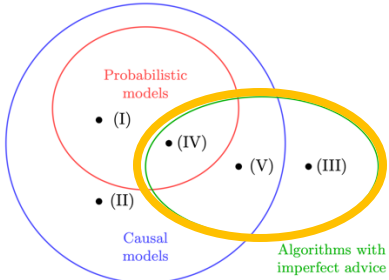
$$\widehat{M}$$



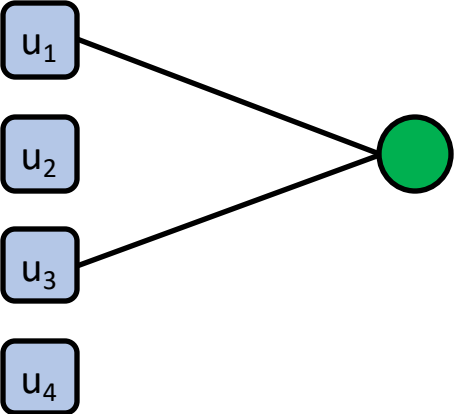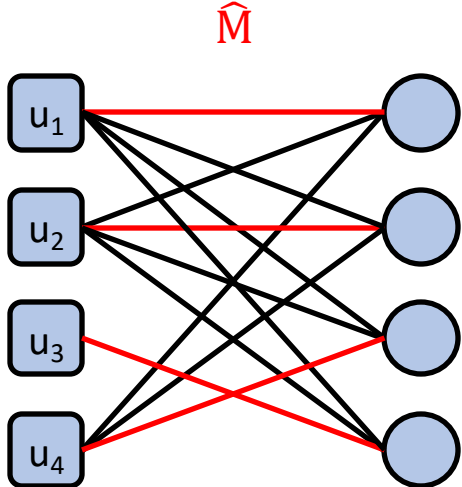| Type | $c^*$ | $\hat{c}$ |
|------|-------|-----------|
| $\{u_1, u_2, u_4\}$ | 2 | 3 |
| $\{u_1, u_3\}$ | 1 | 0 |
| $\{u_2, u_3\}$ | 1 | 0 |
| $\{u_1, u_2, u_3\}$ | 0 | 1 |

# A glimpse of [CGLB24]

Insight: "Testing can be cheaper than learning"

- Mimic algorithm: Fix arbitrary maximum matching $\widehat{M}$ defined by $\hat{c}$ and try to follow it as much as possible. If unable, leave unmatched

$\widehat{M}$

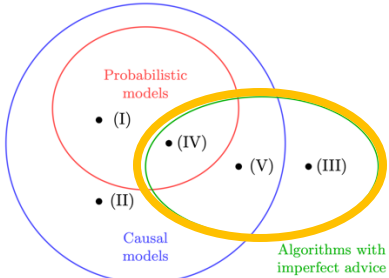| Type | $c^*$ | $\hat{c}$ |
|---|---|---|
| $\{u_1, u_2, u_4\}$ | 2 | 3 |
| $\{u_1, u_3\}$ | 1 | 0 |
| $\{u_2, u_3\}$ | 1 | 0 |
| $\{u_1, u_2, u_3\}$ | 0 | 1 |

# A glimpse of [C̲GLB24]



Insight: "Testing can be cheaper than learning"

- Mimic algorithm: Fix arbitrary maximum matching $\widehat{M}$ defined by $\hat{c}$ and try to follow it as much as possible. If unable, leave unmatched



$\widehat{M}$

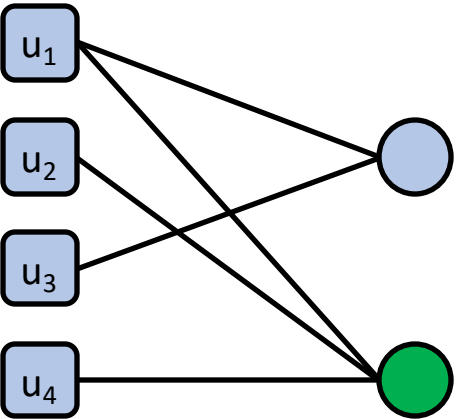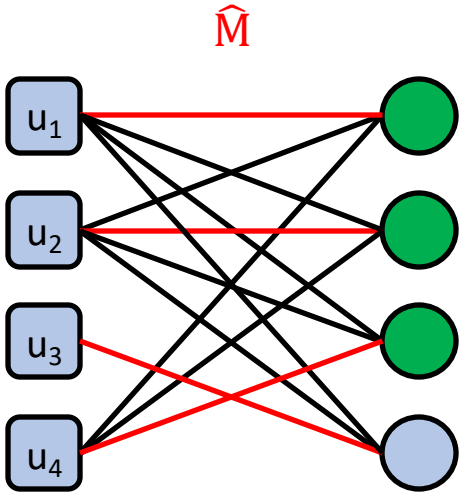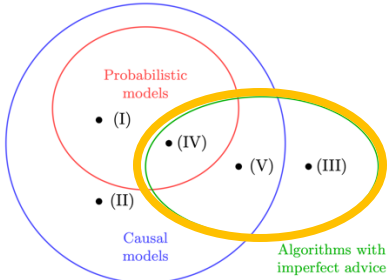| Type | $c^*$ | $\hat{c}$ |
|------|-------|-----------|
| $\{u_1, u_2, u_4\}$ | 2 | 3 |
| $\{u_1, u_3\}$ | 1 | 0 |
| $\{u_2, u_3\}$ | 1 | 0 |
| $\{u_1, u_2, u_3\}$ | 0 | 1 |

# A glimpse of [C̲GLB24]



Insight: "Testing can be cheaper than learning"

- Mimic algorithm: Fix arbitrary maximum matching $\widehat{M}$ defined by $\hat{c}$ and try to follow it as much as possible. If unable, leave unmatched

$$\widehat{M}$$



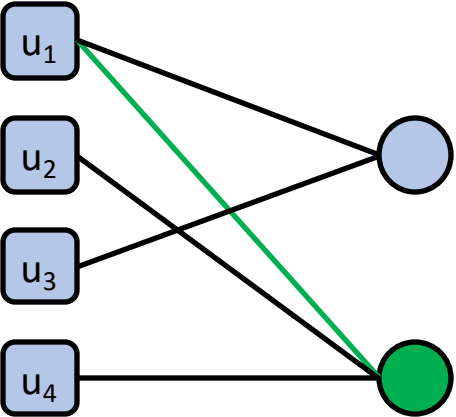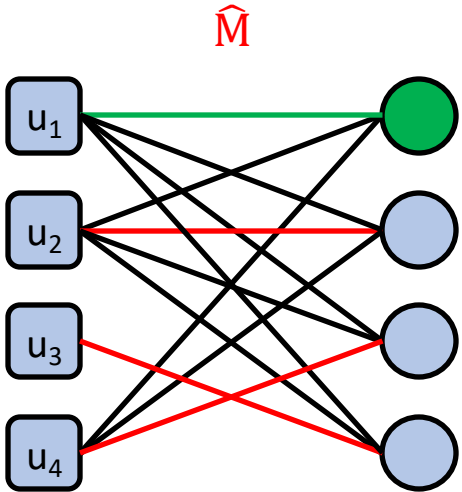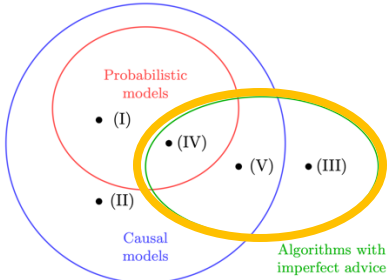| Type | $c^*$ | $\hat{c}$ |
|---|---|---|
| $\{u_1, u_2, u_4\}$ | 2 | 3 |
| $\{u_1, u_3\}$ | 1 | 0 |
| $\{u_2, u_3\}$ | 1 | 0 |
| $\{u_1, u_2, u_3\}$ | 0 | 1 |

# A glimpse of [C̲GLB24]



Insight: "Testing can be cheaper than learning"

- Mimic algorithm: Fix arbitrary maximum matching $\widehat{M}$ defined by $\hat{c}$ and try to follow it as much as possible. If unable, leave unmatched

$\widehat{M}$



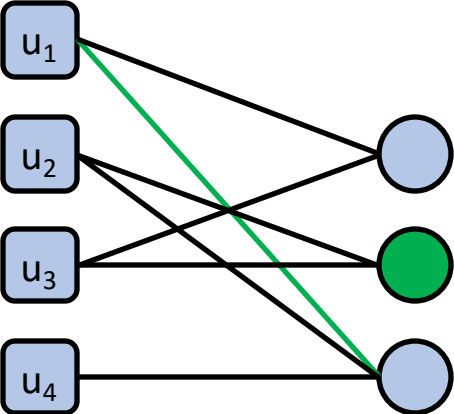| Type | $c^*$ | $\hat{c}$ |
|---|---|---|
| $\{u_1, u_2, u_4\}$ | 2 | ~~3~~ 2 |
| $\{u_1, u_3\}$ | 1 | 0 |
| $\{u_2, u_3\}$ | 1 | 0 |
| $\{u_1, u_2, u_3\}$ | 0 | 1 |

# A glimpse of [C̲GLB24]



**Insight: "Testing can be cheaper than learning"**

- Mimic algorithm: Fix arbitrary maximum matching $\widehat{M}$ defined by $\hat{c}$ and try to follow it as much as possible. If unable, leave unmatched



$\widehat{M}$

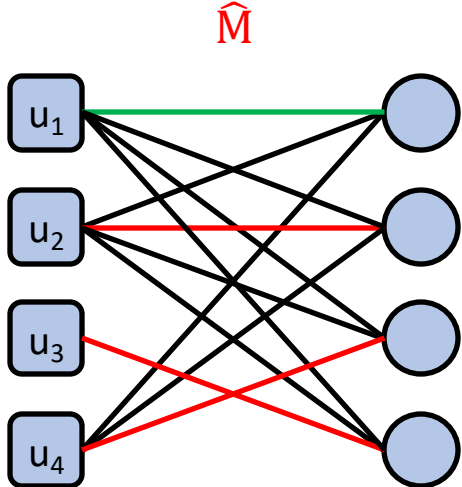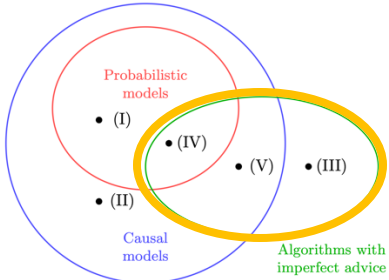| Type | $c^*$ | $\hat{c}$ |
|---|---|---|
| $\{u_1, u_2, u_4\}$ | 2 | ~~3~~ 2 |
| $\{u_1, u_3\}$ | 1 | 0 |
| $\{u_2, u_3\}$ | 1 | 0 |
| $\{u_1, u_2, u_3\}$ | 0 | 1 |

# A glimpse of [CGLB24]



Insight: "Testing can be cheaper than learning"

- Mimic algorithm: Fix arbitrary maximum matching $\widehat{M}$ defined by $\hat{c}$ and try to follow it as much as possible. If unable, leave unmatched

$$\widehat{M}$$



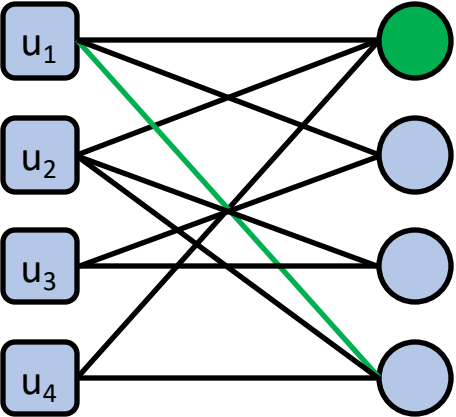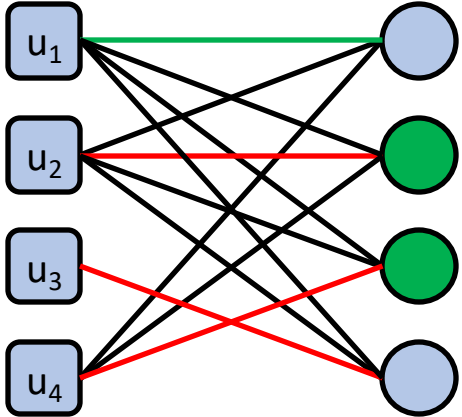| Type | $c^*$ | $\hat{c}$ |
|---|---|---|
| $\{u_1, u_2, u_4\}$ | 2 | ~~3~~ 2 |
| $\{u_1, u_3\}$ | 1 | 0 |
| $\{u_2, u_3\}$ | 1 | 0 |
| $\{u_1, u_2, u_3\}$ | 0 | 1 |

# A glimpse of [C̲GLB24]
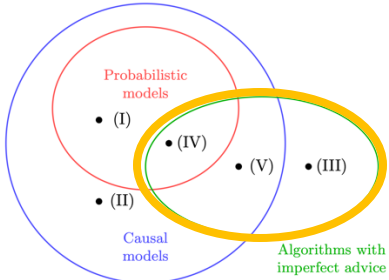
Insight: "Testing can be cheaper than learning"

- Mimic algorithm: Fix arbitrary maximum matching $\widehat{M}$ defined by $\hat{c}$ and try to follow it as much as possible. If unable, leave unmatched

$\widehat{M}$



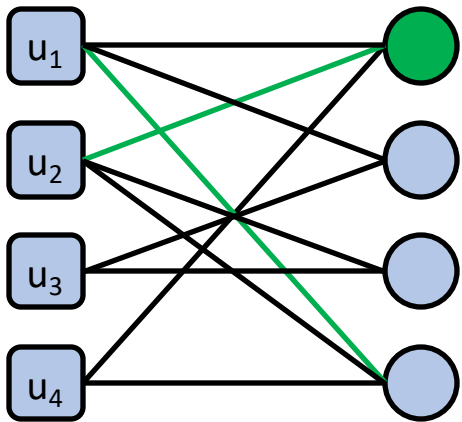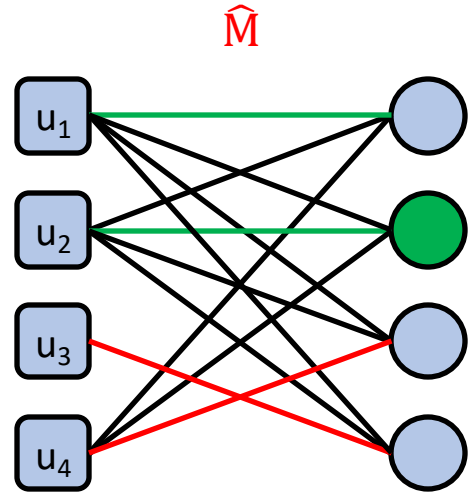| Type | $c^*$ | $\hat{c}$ |
|------|-------|-----------|
| $\{u_1, u_2, u_4\}$ | 2 | ~~3~~ ~~2~~ 1 |
| $\{u_1, u_3\}$ | 1 | 0 |
| $\{u_2, u_3\}$ | 1 | 0 |
| $\{u_1, u_2, u_3\}$ | 0 | 1 |

# A glimpse of [CGLB24]

Insight: "Testing can be cheaper than learning"

- Mimic algorithm: Fix arbitrary maximum matching $\widehat{M}$ defined by $\hat{c}$ and try to follow it as much as possible. If unable, leave unmatched

| Type | $c^*$ | $\hat{c}$ |
|------|-------|-----------|
| $\{u_1, u_2, u_4\}$ | 2 | 3 |
| $\{u_1, u_3\}$ | 1 | 0 |
| $\{u_2, u_3\}$ | 1 | 0 |
| $\{u_1, u_2, u_3\}$ | 0 | 1 |

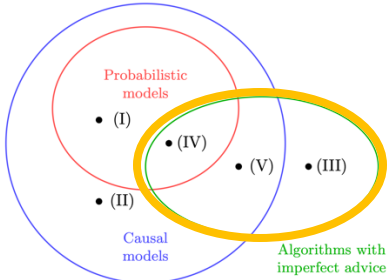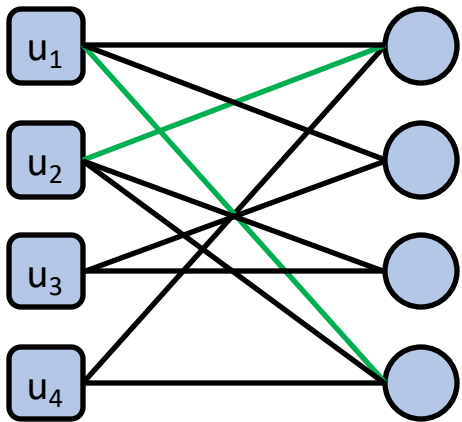Produced matching size

= 2

# A glimpse of [CGLB24]



Insight: "Testing can be cheaper than learning"

- Mimic algorithm: Fix arbitrary maximum matching $\widehat{M}$ defined by $\hat{c}$ and try to follow it as much as possible. If unable, leave unmatched



| Type | $c^*$ | $\hat{c}$ |
|---|---|---|
| $\{u_1, u_2, u_4\}$ | 2 | 3 |
| $\{u_1, u_3\}$ | 1 | 0 |
| $\{u_2, u_3\}$ | 1 | 0 |
| $\{u_1, u_2, u_3\}$ | 0 | 1 |

Produced matching size

= 2

$L_1(c^*, \hat{c})$
$= |3 - 2| + |0 - 1|$
$\ + |0 - 1| + |1 - 0| + 0 \dots$
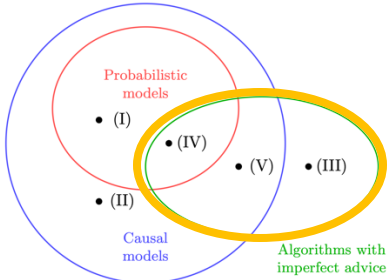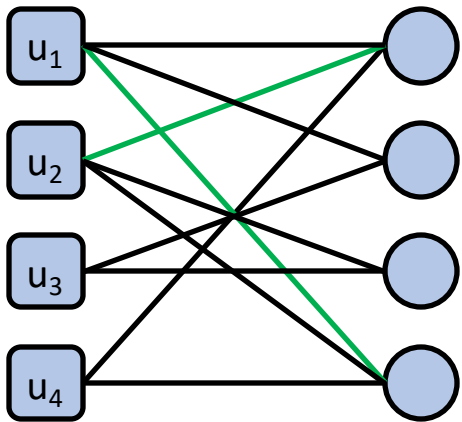$= 4$

# A glimpse of [CGLB24]



Insight: "Testing can be cheaper than learning"

- Mimic algorithm: Fix arbitrary maximum matching $\widehat{M}$ defined by $\hat{c}$ and try to follow it as much as possible. If unable, leave unmatched



| Type | $c^*$ | $\hat{c}$ |
|---|---|---|
| $\{u_1, u_2, u_4\}$ | 2 | 3 |
| $\{u_1, u_3\}$ | 1 | 0 |
| $\{u_2, u_3\}$ | 1 | 0 |
| $\{u_1, u_2, u_3\}$ | 0 | 1 |

Produced matching size
$$= 2 = \left|\widehat{M}\right| - \frac{L_1(c^*, \hat{c})}{2}$$

Error is "double counted" in $L_1$

$$L_1(c^*, \hat{c})$$
$$= |3 - 2| + |0 - 1|$$
$$+ |0 - 1| + |1 - 0| + 0 \ldots$$
$$= 4$$

# A glimpse of [CGLB24]



Insight: "Testing can be cheaper than learning"

- Mimic algorithm: Fix arbitrary maximum matching $\widehat{M}$ defined by $\hat{c}$ and try to follow it as much as possible. If unable, leave unmatched
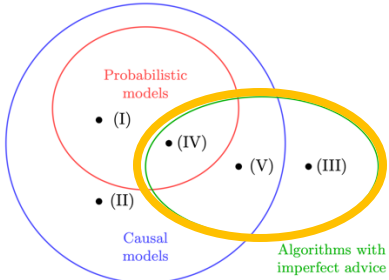
- Analysis: $0 \leq L_1(c^*, \hat{c}) \leq 2n$ measures how close $\hat{c}$ is to $c^*$

  - By blindly following advice, Mimic gets a matching of size $\left|\widehat{M}\right| - \frac{L_1(c^*, \hat{c})}{2}$

  - Mimic beats an advice-free Baseline whenever $\left|\widehat{M}\right| - \frac{L_1(c^*, \hat{c})}{2} > \beta \cdot n$

# A glimpse of [CGLB24]



Insight: "Testing can be cheaper than learning"

- Mimic algorithm: Fix arbitrary maximum matching $\widehat{M}$ defined by $\hat{c}$ and try to follow it as much as possible. If unable, leave unmatched
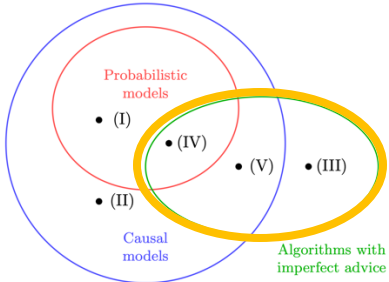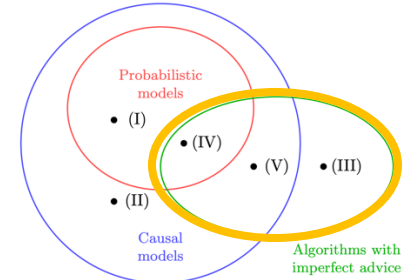
- Analysis: $0 \leq L_1(c^*, \hat{c}) \leq 2n$ measures how close $\hat{c}$ is to $c^*$

  - By blindly following advice, Mimic gets a matching of size $\left|\widehat{M}\right| - \frac{L_1(c^*, \hat{c})}{2}$

  - Mimic beats an advice-free Baseline whenever $\left|\widehat{M}\right| - \frac{L_1(c^*, \hat{c})}{2} > \beta \cdot n$

- Idea: Use Mimic when $L_1(c^*, \hat{c})$ low; otherwise use Baseline

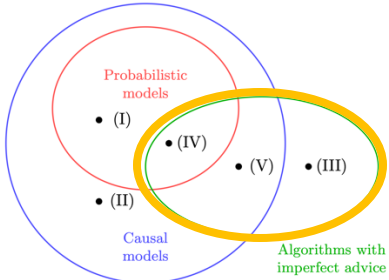- Problem: We don't know $c^*$, so cannot evaluate $L_1(c^*, \hat{c})$

# A glimpse of [CGLB24]



**Insight: "Testing can be cheaper than learning"**

- Use sublinear property testing to estimate $L_1(c^*, \hat{c})$ ← *Random arrival ordering* $\equiv$ i.i.d. samples

  - Define $p = \frac{c^*}{n}$ and $q = \frac{\hat{c}}{n}$ as distributions over the $2^U$ types

  - [VV11, JHW18]: Can estimate $L_1(p, q)$ "well" using o(n) i.i.d. samples

    - Some adjustments needed to our problem setting, but it can be done

[VV11] Gregory Valiant and Paul Valiant. *The power of linear estimators*. Foundations of Computer Science (FOCS), 2011.
[JHW18] Jiantao Jiao, Yanjun Han, and Tsachy Weissman. *Minimax estimation of the L₁ distance*. IEEE Transactions on Information Theory, 2018.

# A glimpse of [CGLB24]

Lower bound on achieved competitive ratio (w.p. $\geq 1 - \delta$)

per than learning"

ting to estimate $L_1(c^*, \hat{c})$ ← Random arrival ordering $\equiv$ i.i.d. samples

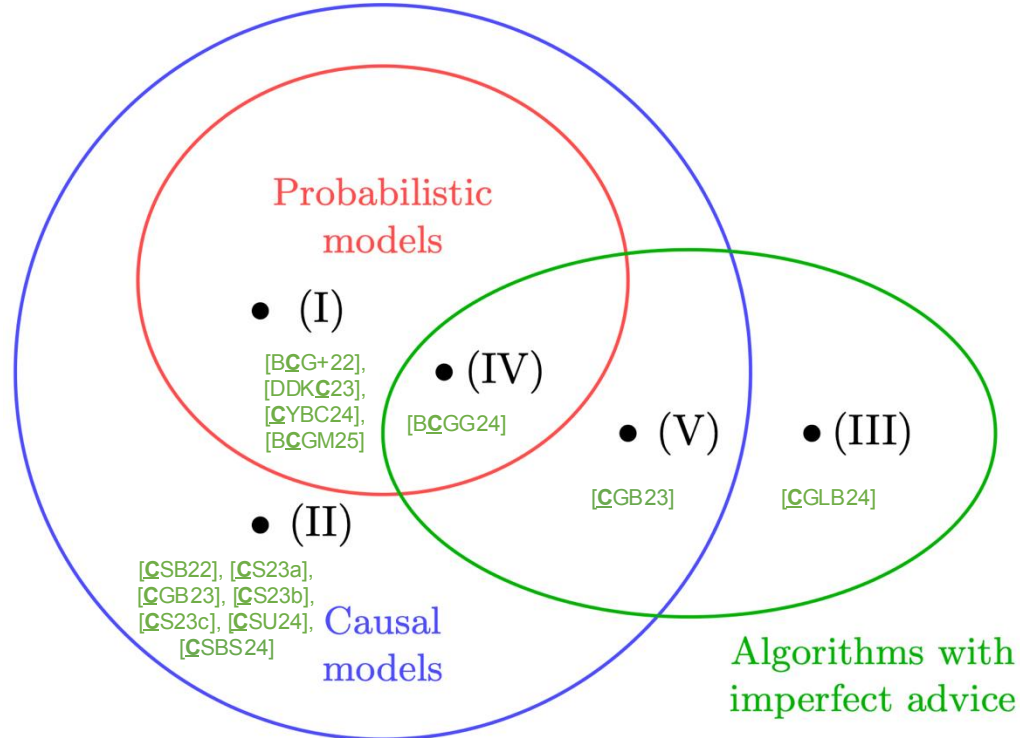$\beta \cdot (1 - o_n(1))$ distributions over the $2^U$ types

te $L_1(p, q)$ "well" using o(n) i.i.d. samples

o our problem setting, but it can be done

$2\left(\frac{\hat{n}}{n} - \beta\right) - 2\varepsilon$ ... $1$ ... $L_1(p^*, q)$

- **TestAndMatch**: Use Mimic or Baseline depending on $\hat{L}_1(c^*, \hat{c})$

  - Achieve comp. ratio at least $1 - \frac{L_1(c^*, \hat{c})}{2n} \geq \beta$, when $\hat{L}_1(c^*, \hat{c})$ "small"

  - Achieve comp. ratio at least $\beta \cdot (1 - o(1))$, when $\hat{L}_1(c^*, \hat{c})$ "large"

  - i.e., TestAndMatch is 1-consistent and $\beta \cdot (1 - o(1))$-robust

# Main themes explored in my PhD thesis



Probabilistic
models

● (I)

[B**C**G+22],
[DDK**C**23],
[**C**YBC24],
[B**C**GM25]

● (IV)

[B**C**GG24]

● (V)

[**C**GB23]

● (III)

[**C**GLB24]

● (II)

[**C**SB22], [**C**S23a],
[**C**GB23], [**C**S23b],
[**C**S23c], [**C**SU24],
[**C**SBS24]

Causal
models

Algorithms with
imperfect advice

# Research vision:
# Principled algorithms with real-world impact



Causality-aware AI/ML methods

Algorithms with imperfect advice

Applying insights to real-world problems

**Theory**

**Practice**

# Research vision:
# Principled algorithms with real-world impact

Causality-aware
AI/ML methods

**Theory**

- Beyond simple statistical or association relationships
- Especially important for systems that act on the environment and have impact on real-world decisions

# Research vision:
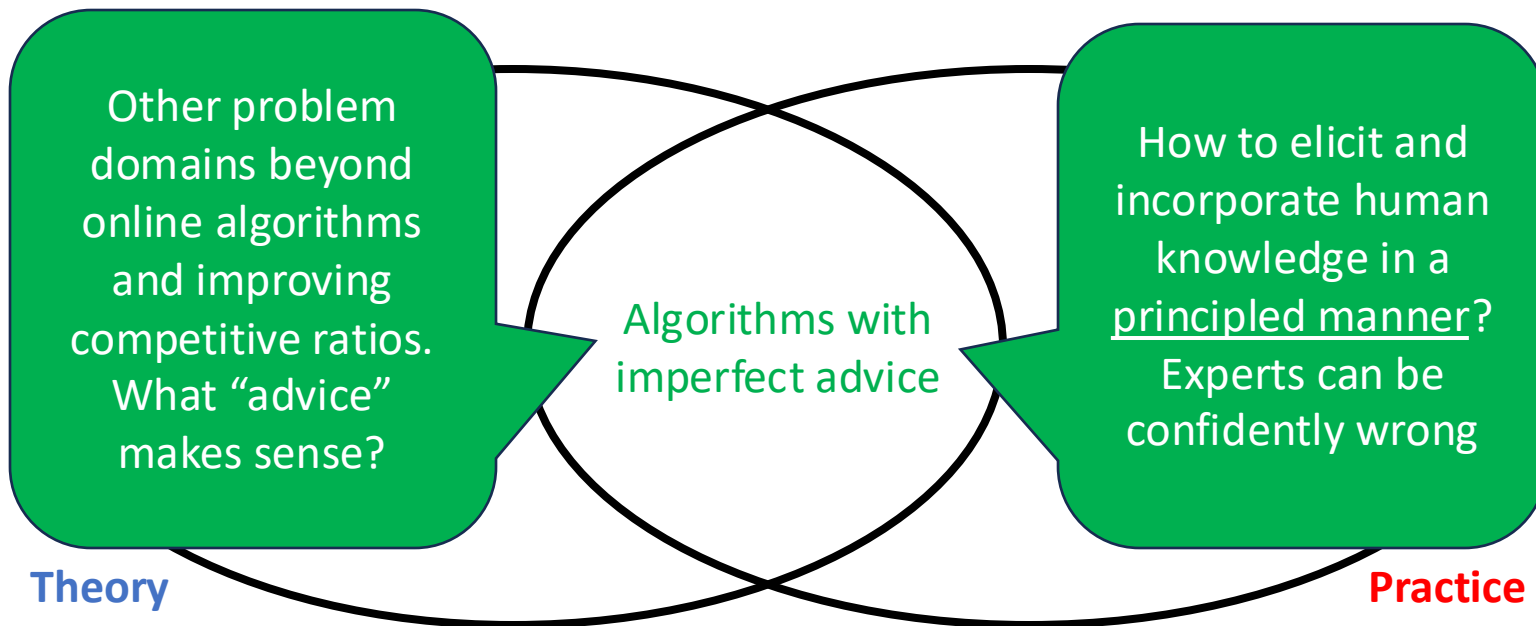# Principled algorithms with real-world impact



- Translate knowledge to benefit society
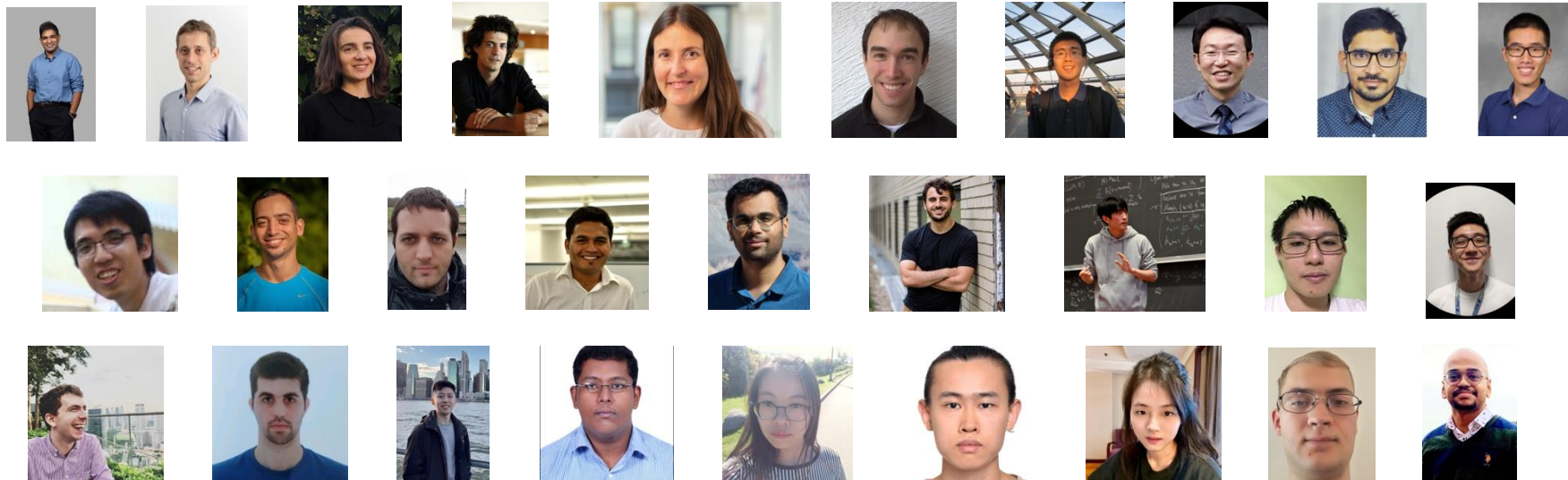- Model and solve real-world problems in a principled manner

Applying insights to real-world problems

**Practice**

# Research vision:
# Principled algorithms with real-world impact

# Thank you to all my amazing collaborators during my PhD journey!



**Thank you for your kind attention!**