# (On the Power of) Knowledge Compilation in Causal Inference

Alessandro Antonucci (alessandro@idsia.ch)
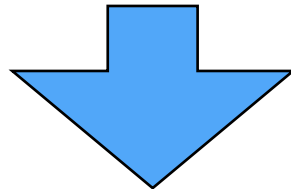
IDSIA USI-SUPSI

*Workshop on Probabilistic Circuits and Logic - Simons Institute, Berkeley - Oct 17, 2023*

IDSIA

SUPSI

# What is this Talk about

"*Knowledge compilation has been successfully used to solve beyond NP problems, including some PP-complete and NP$^{PP}$-complete problems for Bayesian networks.*"

*Solving PP$^{PP}$-complete problems using knowledge compilation*, Otzok, Choi, and Darwiche (KR, 2016)

6th Workshop on Tractable Probabilistic Modeling
Building Bridges
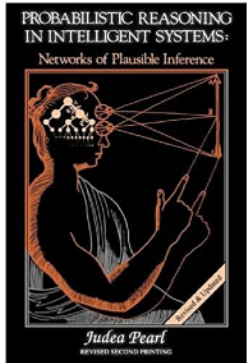
[Submitted on 5 Oct 2023]

## Tractable Bounding of Counterfactual Queries by Knowledge Compilation

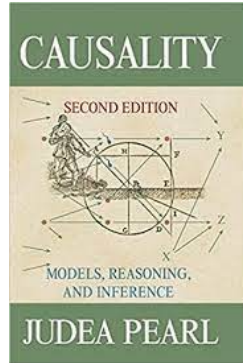David Huber, Yizuo Chen, Alessandro Antonucci, Adnan Darwiche, Marco Zaffalon
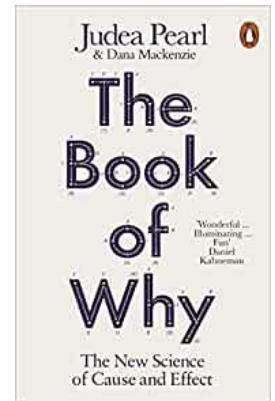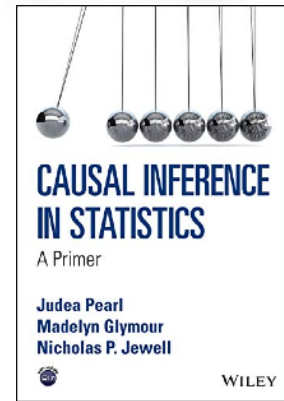
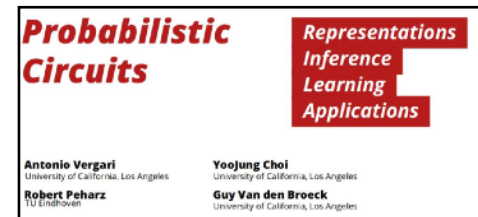# From Pearl to Pearl

**Pearl**

Bayesian Nets
( ~1988 )

Do Calculus
( ~2000 )

**This Talk**

Structural Causal Models
( ~2016 )

**Darwiche**

Knowledge Compilation
( ~2000 )

"TPM" Community ( ~2020 )

# ( Science > ) AI > Deep Learning



Andrej Karpathy blog
About
The Unreasonable Effectiveness of Recurrent Neural Networks
May 21, 2015

RESEARCH ARTICLE | BIOLOGICAL SCIENCES

The unreasonable effectiveness of deep learning in artificial intelligence

Terrence J. Sejnowski | Authors Info & Affiliations
Edited by David L. Donoho, Stanford University, Stanford, CA, and approved November 22, 2019 (received for review September 17, 2019)

January 28, 2020 | 117 (48) 30033-30038 | https://doi.org/10.1073/pnas.1907373117

*"Deep learning has instead given us machines with truly **impressive abilities but no intelligence**.*

Pearl

*The difference is profound and lies in the **absence of a model of reality**."*



COMMUNICATIONS of the ACM

Human-Level Intelligence or Animal-Like Abilities?

Computing within Limits
Transient Electronics Take Shape
Q&A with Dina Katabi
Formally Verified Software in the Real World

Darwiche

Judea Pearl & Dana Mackenzie
The Book of Why
The New Science of Cause and Effect



Make it RAIN!!!

Correlation ≠ Causation

OBSERVER

www.thegraphicrecorder.com

# Pearl's Ladder of Causation and the Need for a Causal AI

## 3-LEVEL HIERARCHY

(Causal) AI?

RL

ML/DL



IMAGINING

DOING

SEEING

3. COUNTERFACTUALS
ACTIVITY:     Imagining, Retrospection, Understanding
QUESTIONS:    *What if I had done . . . ? Why?*
              (Was it X that caused Y? What if X had not
              occurred? What if I had acted differently?)
EXAMPLES:     Was it the aspirin that stopped my headache?
              Would Kennedy be alive if Oswald had not
              killed him? What if I had not smoked the last 2 years?

2. INTERVENTION
ACTIVITY:     Doing, Intervening
QUESTIONS:    *What if I do . . . ? How?*
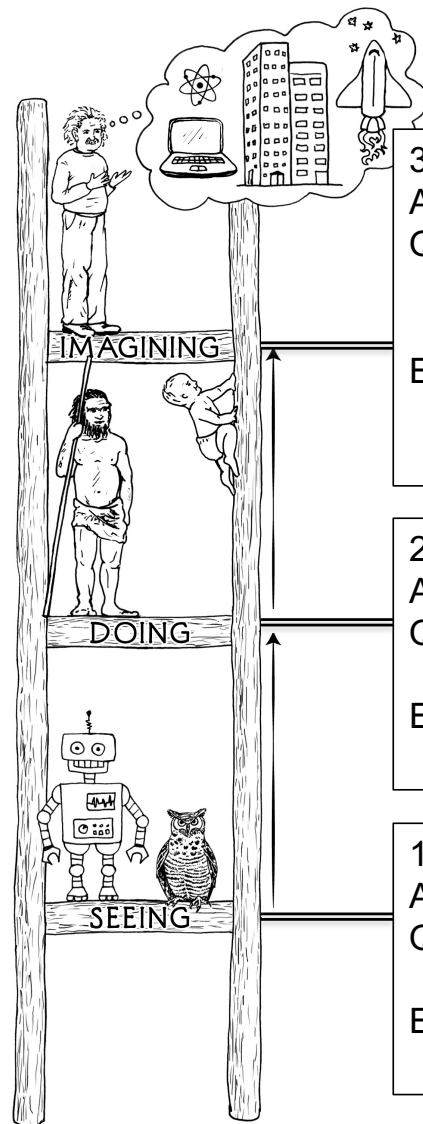              (What would Y be if I do X?)
EXAMPLES:     If I take aspirin, will my headache be cured?
              What if we ban cigarettes?

1. ASSOCIATION
ACTIVITY:     Seeing, Observing
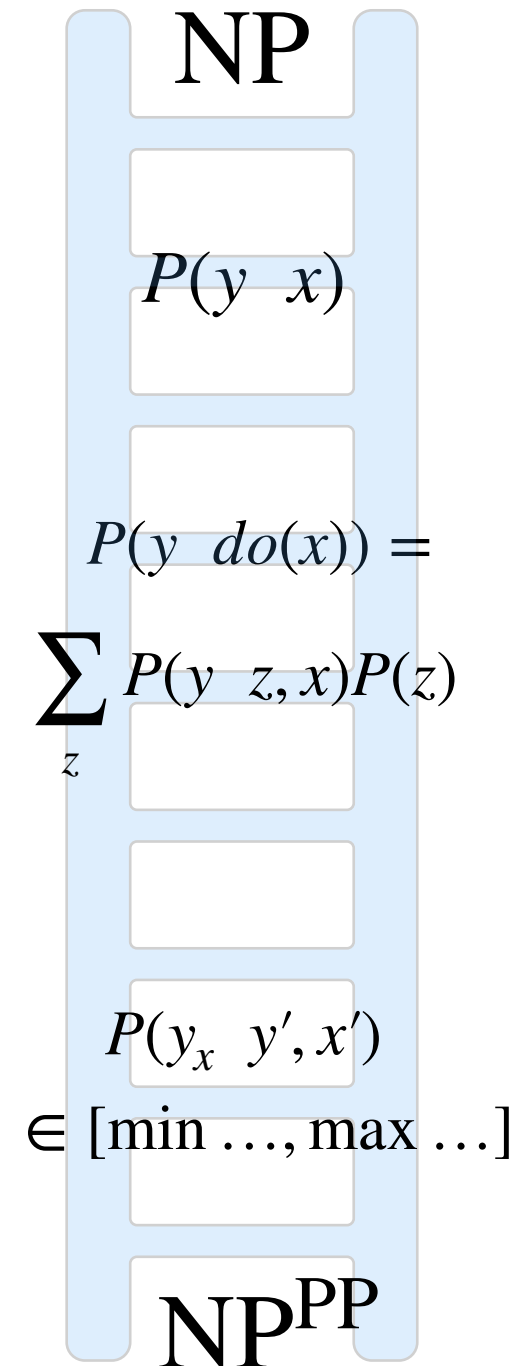QUESTIONS:    *What if I see . . . ?*
              (How would seeing X change my belief in Y?)
EXAMPLES:     What does a symptom tell me about a disease?
              What does a survey tell us about the election results?

Source: The Book of Why, Pearl & Mc Kenzie

# A Ladder for (PGM) Inference?

- Answering an **observational** query?
  Single PGM query in the empirical model

- (Identifiable) **interventional** query?
  - **Do-calculus** and queries by auxiliary PGM inferences in the empirical model
  - Single EM on the SCM with latent variables + PGM inference (Dechter, 2023)

- **Counterfactual** queries suffer partial identifiability (**bounds** only)
  - Credal nets (Zaffalon & Antonucci, 2020)
  - Multiple EM runs (Zaffalon & Antonucci, 2021)
  - Sampling (Bareinboim, 2022)
  - Polynomial programs (Shpitser, 2023)
  - Multiple EM + KG (this talk)

$$\text{NP}$$

$$P(y \mid x)$$

$$P(y \mid do(x)) =$$

$$\sum_z P(y \mid z, x) P(z)$$

$$P(y_x \mid y', x')$$

$$\in [\min \ldots, \max \ldots]$$

$$\text{NP}^{\text{PP}}$$

# Structural Causal Models (Univariate)

- Manifest **endogenous** variable $X$

- Observations $\mathcal{D}$ available

- From $\mathcal{D}$ statistical learning of $P(X)$

$$X$$

Boolean $X$
$$P(X = 0) = p$$

# Structural Causal Models (Univariate)

- Manifest **endogenous** variable $X$

- Observations $\mathscr{D}$ available

- From $\mathscr{D}$ statistical learning of $P(X)$

- A latent **exogenous** variable $U$

- $U$ determines $X$ (**structural equation** $f_X$)

- $P(U)$ induces (a single) $P(X)$

$$P(x) = \sum_x P(x \mid u)P(u) = \sum_u \delta_{f(u),x}P(u)$$

$U \in \{0,1,2,3\}$

U

$f_X$

$f_X(U = 0) = 0$
$f_X(U = 1) = 0$
$f_X(U = 2) = 1$
$f_X(U = 3) = 1$

X

Boolean $X$

$P(X = 0) = p$

# Structural Causal Models (Univariate)

- Manifest **endogenous** variable $X$

$$P(U) = \left[\frac{p}{2}, \frac{p}{2}, \frac{1-p}{2}, \frac{1-p}{2}\right]$$

- Observations $\mathcal{D}$ available

$$U \in \{0,1,2,3\}$$

- From $\mathcal{D}$ statistical learning of $P(X)$

- A latent **exogenous** variable $U$

- $U$ determines $X$ (**structural equation** $f_X$)

$f_X(U = 0) = 0$

$f_X(U = 1) = 0$

$f_X(U = 2) = 1$

$f_X(U = 3) = 1$

- $P(U)$ induces (a single) $P(X)$

$$P(x) = \sum_x P(x \mid u)P(u) = \sum_u \delta_{f(u),x} P(u)$$

- $P(X)$ to $P(U)$? Multiple consistent $P(U)$'s

Boolean $X$

$P(X = 0) = p$

- Bounds? Query has different values for the different consistent $P(U)$!

# Structural Causal Models

- $\mathbf{X} := (X_1, \ldots, X_n)$ (endogenous variables)
- $\mathbf{U} := (U_1, \ldots, U_m)$ (exogenous variables)
- Directed graph $\mathscr{G}$ assumed to be semi-Markovian = root in $\mathbf{U}$, non-root in $\mathbf{X}$

- Equation $X = f_X(\mathrm{Pa}_X)$ for each $X \in \mathbf{X}$
  - Exogenous states $\mathbf{U} = \mathbf{u}$ determine endogenous states $\mathbf{X} = \mathbf{x}$
- Marginal $P(U)$ for each $U \in \mathbf{U}$
  - Exogenous distribution distribution $P(\mathbf{U})$ induces endogenous distribution $P(\mathbf{X})$

$$U = \{X, Y\}, \quad V = \{Z\}, \quad F = \{f_Z\}$$
$$f_Z : Z = 2X + 3Y$$

# SCMs as BNs?

- **An SCM is a BN** with CPTs $P(X \mid \text{Pa}_X) = \delta_{X, f_X(\text{Pa}_X)}$

$$P(\mathbf{x}, \mathbf{u}) = \prod_{U \in \mathbf{U}} P(u) \prod_{X \in \mathbf{X}} \delta_{f_X(\text{pa})_X, x}$$

- We need:
  – Causal Graph (= Exogenous Confounders)
  – Structural Equations (= Endogenous CPTs)
  – Exogenous Marginals

**FSCM = Fully Specified**

- Often we only have:
  – Causal Graph
  – Endogenous Data
  – Structural Equations? "Canonical" specification

**PSCM = Partially Specified**

# Canonical Specification of Structural Equations

- Structural equations from $\mathcal{G}$ ?

- $y = f(x, u)$? Canonical? $U$ indexing all deterministic mechanisms btw X and Y

- With Boolean parent & child?

- $U = 4$

- In general, exponential size:

$$U = Y^{\prod_{X \in \text{Pa}_Y} X}$$

- Even larger cardinality if Y has more than an exogenous parent

latent



$$y = f(x, u)$$

ex. **disease** and test **outcome**

$$P(Y \mid X, U)$$

| | X=0 | X=1 | X=0 | X=1 | X=0 | X=1 | X=0 | X=1 |
|---|---|---|---|---|---|---|---|---|
| Y=0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 0 |
| Y=1 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 1 |
| | U=0 | | U=1 | | U=2 | | U=3 | |

$$Y = 0 \qquad Y = X \qquad Y = \neg X \qquad Y = 1$$
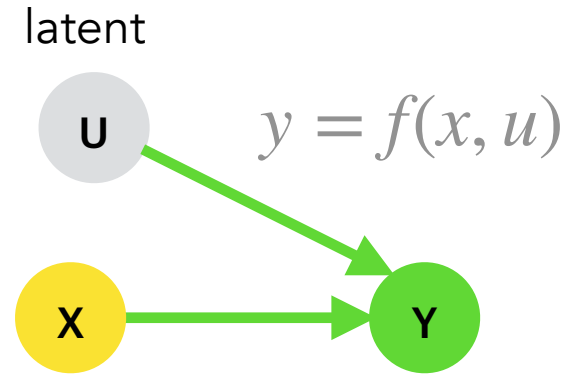
# Canonical Specification of Structural Equations

- Structural equations from $\mathcal{G}$ ?

- $y = f(x, u)$? Canonical? $U$ indexing all deterministic mechanisms btw X and Y

- With Boolean parent & child?

- $\ \ U\ = 4$

- In general, exponential size:

$$U\ =\ Y^{\ \prod_{X \in \text{Pa}_Y} X}$$

- Even larger cardinality if Y has more than an exogenous parent

- Non-canonical? Domain knowledge (ex. $Y = 1$ and $Y = \neg X$ impossible)

latent



$y = f(x, u)$

ex. **disease** and test **outcome**

$P(Y\ X, U)$

| | X=0 | X=1 | X=0 | X=1 | X=0 | X=1 | X=0 | X=1 |
|-----|-----|-----|-----|-----|-----|-----|-----|-----|
| Y=0 | 1 | 1 | 1 | 0 | | | | |
| Y=1 | 0 | 0 | 0 | 1 | | | | |
| | U=0 | | U=1 | | U=2 | | U=3 | |

$Y = 0 \qquad Y = X \qquad Y = \neg X \qquad Y = 1$

# Inference in FSCMs

- BN inference is $O(2^{\text{treewidth}})$, faster with:
  - **c**ontext-**s**pecific **i**ndependence
  - **det**erminism
- FSCM = BN + determinism in CPTs

  - Compilation to tractable circuits with FSCMs of high tw (>100)
  - Causal treewidth $\leq$ treewidth inference $O(2^{\text{causal treewidth}})$

| Local Structure Encoded | Pathfinder | Water | Munin4 |
|---|---|---|---|
| None | 981,178 | 13,777,166 | 116,136,985 |
| Det + CSI | 42,810 (4%) | 134,140 (1%) | 5,762,690 (5%) |
| Det | 130,380 (13%) | 138,501 (1%) | 9,997,267 (9%) |
| CSI | 200,787 (20%) | 11,111,104 (81%) | 17,612,036 (15%) |

- Operational characterisation (Darwiche, 2022)
- Counterfactuals? ctw x (# of worlds)
- Standard compilers (ex. ACE) not specialized to FSCMs

# Inference in PSCMs

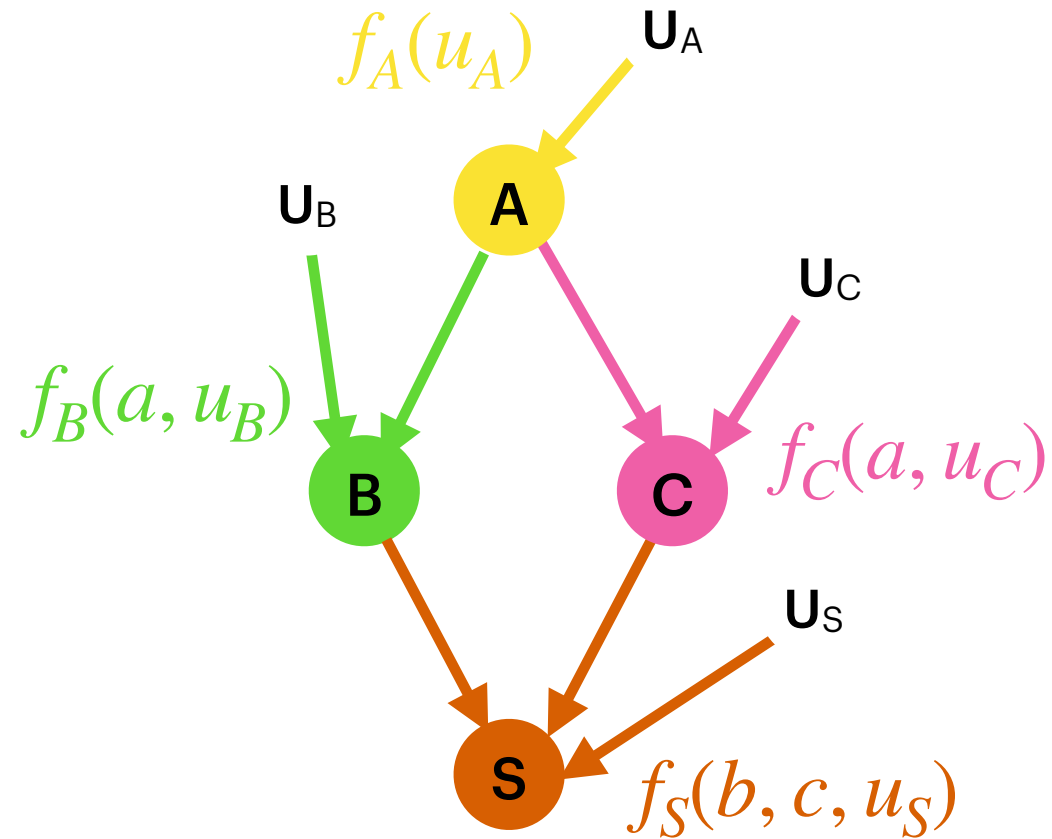- More challenging than FSCM inference

- Identifiable queries?

    – Do-calculus = inference in the empirical BN

- Non-identifiable?

    – Bound computation

    – Equivalent to inference in a **credal net**
      (i.e., bounds wrt iterated BN inference)

    – NP$^{PP}$ task (Zaffalon and Antonucci, 2023)

- PSCM = Collection of compatible FSCMs

- Let's write the compatibility constraints!

# Credal Net Mapping

- Find the exogenous marginals?

  $$P(U_A)P(U_B)P(U_C)P(U_S)$$

- **Endogenous** (= with $\mathcal{D}$) **consistency**

- This induces global non-linear (so-called Verma) constraints

- Let's make the constraints local and linear by marginalisation and conditioning



$$\sum_{u_A, u_B, u_C, u_D} \left[ p(u_A) \cdot \delta_{a, f_A(u_A)} \cdot p(u_B) \cdot \delta_{b, f_B(a, u_B)} \cdot p(u_C) \cdot \delta_{c, f_C(a, u_c)} \cdot p(u_S) \cdot \delta_{s, f_S(b, c, u_S)} \right] = \tilde{p}(a, b, c, s)$$

Unknown    Unknown    Unknown    Unknown    Empirical, known

# Credal Net Mapping (con't)

$$\sum_{u_A,u_B,u_C,u_D} \left[ p(u_A) \cdot \delta_{a,f_A(u_A)} \cdot p(u_B) \cdot \delta_{b,f_B(a,u_B)} \cdot p(u_C) \cdot \delta_{c,f_C(a,u_C)} \cdot p(u_S) \cdot \delta_{s,f_S(b,c,u_S)} \right] = \tilde{p}(a,b,c,s)$$

$f_A(u_A)$ $\quad$ $U_A$

$U_B$ $\quad$ **A**

$U_C$

$f_B(a,u_B)$

$f_C(a,u_C)$

**B** $\qquad$ **C**

$U_S$

$$P(a) = \sum_{u_A} P(a \mid u_A) \cdot P(u_A)$$

$$P(b \mid a) = \sum_{u_B} P(b \mid a, u_B) \cdot P(u_B)$$

$$P(c \mid a) = \sum_{u_C} P(c \mid a, u_C) \cdot P(u_C)$$

$$P(s \mid b, c) = \sum_{u_S} P(s \mid b, c, u_S) \cdot P(u_S)$$

**S**

$f_S(b,c,u_S)$

- Linear constraints on marginal exogenous probabilities leading to the (credal) set specification $K(U_A)$, $K(U_B)$, $K(U_C)$, $K(U_S)$
- Structural equations (= endogenous CPTS) remain unaffected

# Causal Inference by Credal Nets



$K(U_A)$

$K(U_B)$

A

$K(U_C)$

B

C

$K(U_S)$

S

$$P(B \mid do(\overline{a})) \in [\underline{P}'(B \mid \overline{a}), \overline{P}'(B \mid \overline{a})]$$

Interventional query

- Identifiable? $\underline{P} = \overline{P}$

# Causal Inference by Credal Nets

$$K(U_A)$$

$$K(U_B)$$

$$K(U_C)$$

$$P(S_b \ \overline{b}) \in [\underline{P}(S \ b, \overline{b}'), \overline{P}(S \ b, \overline{b}')]$$

A

Counterfactual query

B

C

B

$$K(U_S)$$

B'

S

S

$$P(B \ \text{do}(\overline{a})) \in [\underline{P}'(B \ \overline{a}), \overline{P}'(B \ \overline{a})]$$

Interventional query

- Identifiable? $\underline{P} = \overline{P}$

# Causal EM (Zaffalon & Antonucci, 2021)

- CN mapping suffers in models with multiple exogenous parents

- Exogenous variables are always missing (MAR, asystematic, way)

- Expectation Maximisation (Dempster, 1977)

  – Random initialisation of $P(U)$

  – E-step: Missing data completion by expected (fractional) counts

  – M-step: "completed" data to retrain $P(U)$

  – Iterate until convergence

- EM goes to a (local/global) max of log-lik

| U1 | U2 | X1 | X2 | n |
|----|----|----|----|----|
| * | * | 0 | 0 | ... |
| * | * | 0 | 1 | ... |
| * | * | 1 | 0 | ... |
| * | * | 1 | 1 | ... |

1: $t \leftarrow 0$
2: **while** $P(\mathscr{D}|\{\theta_U^{t+1}\}_{U \in U}) \geq P(\mathscr{D}|\{\theta_U^t\}_{U \in U})$ **do**
3:    **for** $U \in U$ **do**
4:       $\theta_U^{t+1} \leftarrow |\mathscr{D}|^{-1}\sum_{v \in \mathscr{D}}\theta_{U|v}^t$
5:       $t \leftarrow t+1$
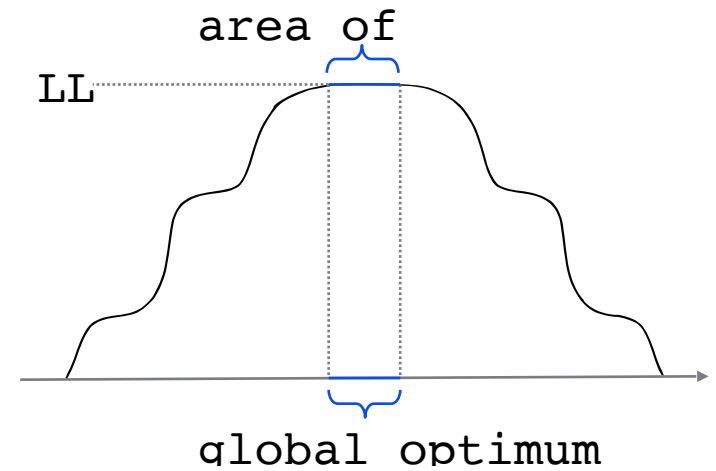6:    **end for**
7: **end while**

# Causal EM: Getting an Inner Approximation of the Bounds

- Causal EM converge to global maximum (that we know) if and only if the corresponding $P(U)$ belongs to credal set $K(U)$

- We sample initialisations, to sample $K(U)$

- For each sample we obtain an inner point

**Theorem 1.** *Let $\mathcal{K}$ denote the set of quantifications for $\{P(U)\}_{U \in U}$ consistent with the following constraint to be satisfied for each $c \in \mathcal{C}$ and each $y^{(c)}$:*

$$(8) \qquad \sum_{\substack{u^{(c)}: f_X(pa_X)=x \\ \forall X \in \mathbf{X}^{(c)}}} \prod_{U \in \mathbf{U}^c} P(u) = \prod_{X \in \mathbf{X}^{(c)}} \hat{P}(x | y_X^{(c)}),$$

*where the values of $u$, $x$ and $y_X^{(c)}$ are those consistent with $u^{(c)}$ and $y^{(c)}$. If $\mathcal{K} \neq \emptyset$, the log-likelihood in Eq. (7) achieves its global maximum if and only if $\{P(U)\}_{U \in U} \in \mathcal{K}$. If $\mathcal{K} = \emptyset$, the marginal log-likelihood in Eq. (7) can only take values strictly lower than the global maximum.*



area of

LL

global optimum

# Causal EM: Getting an Inner Approximation of the Bounds

- Causal EM converge to global maximum (that we $\ldots$) d only if the corresponding $P(U)$ belon $\ldots$

- We sample initialisations $\ldots$

- For each samp $\ldots$

**Theorem 5.** Let $[a^*, b^*]$ denote the exact probability bounds of a causal query. Say that $\rho := \{r_i\}_{i=1}^n$ are the outputs of $n$ EMCC iterations, while $[a, b]$ is the interval induced by $\rho$, i.e., $a := \min_{i=1}^n r_i$ and $b := \max_{i=1}^n r_i$. By construction $a^* \leq a \leq b \leq b^*$. The following inequality holds:

$$P\left(a - \varepsilon L \leq a^* \leq b^* \leq b + \varepsilon L \,\middle|\, \rho\right) = \frac{1 + (1 + 2\varepsilon)^{2-n} - 2(1 + \varepsilon)^{2-n}}{(1 - L^{n-2}) - (n-2)(1 - L)L^{n-2}}, \qquad (13)$$

where $L := (b - a)$ and $\varepsilon := \delta/(2L)$ is the relative error at each extreme of the interval obtained as a function of the absolute allowed error $\delta \in (0, L)$.

**Theo** $\ldots$
follow $\ldots$

(8)

where the $\ldots$ ose consistent with $\mathbf{u}^{(c)}$ and $\mathbf{y}^{(c)}$. If $\mathcal{K} \neq \emptyset$, the
log-likeliho $\ldots$ es its global maximum if and only if $\{P(U)\}_{U \in \mathbf{U}} \in \mathcal{K}$. If
$\mathcal{K} = \emptyset$, the $\ldots$ log-likelihood in Eq. (7) can only take values strictly lower than the
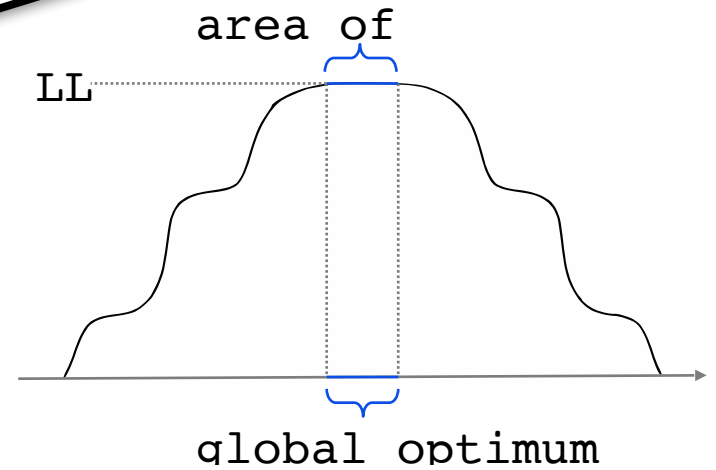global maximum.

area of

LL $\cdots\cdots\cdots$

global optimum

# Causal EM: Getting an Inner Approximation of the Bounds

- Causal EM converge to global maximum (that we ... d only if the corresponding $P(U)$ belon...
- We sample initialisations ...
- For each samp...

**Theorem 5.** Let $[a^*, b^*]$ denote the exact probability bounds of a causal query. Say that $\rho := \{r_i\}_{i=1}^n$ are the outputs of $n$ EMCC iterations, while $[a, b]$ is the interval induced by $\rho$, i.e., $a := \min_{i=1}^n r_i$ and $b := \max_{i=1}^n r_i$. By construction $a^* \le a \le b \le b^*$. The following inequality holds:

$$P\left(a - \varepsilon L \le a^* \le b^* \le b + \varepsilon L \,\middle|\, \rho\right) = \frac{1 + (1+2\varepsilon)^{2-n} - 2(1+\varepsilon)^{2-n}}{(1 - L^{n-2}) - (n-2)(1-L)L^{n-2}}, \quad (13)$$

where $L := (b-a)$ and $\varepsilon := \delta/(2L)$ is the relative error at each extreme of the interval obtained as a function of the absolute allowed error $\delta \in (0, L)$.

area of

LL

20 EM runs to get close to the actual bounds with 95% credibility
For identifiable queries 9 runs to be sure with 99% credibility

# Causal EM (Inferences)

1: $t \leftarrow 0$

2: $\{\theta_U^0\}_{U \in \mathbf{U}} \leftarrow$ random initialisation

3: **while** $P(\mathcal{D}|\{\theta_U^{t+1}\}_{U \in \mathbf{U}}) \geq P(\mathcal{D}|\{\theta_U^t\}_{U \in \mathbf{U}})$ **do**

4:     **for** $U \in \mathbf{U}$ **do**

5:         $\theta_U^{t+1} \leftarrow |\mathcal{D}|^{-1} \sum_{\boldsymbol{x} \in \mathcal{D}} \theta_{U|\boldsymbol{x}}^t$

6:         $t \leftarrow t + 1$

7:     **end for**

8: **end while**

9: **return** $\{\theta_U^{t+1}\}_{U \in \mathbf{U}}$

This is a single run, returning exogenous chances
to be iterated for different random initialisations

# Causal EM (Inferences)

1: $t \leftarrow 0$
2: $\{\theta_U^0\}_{U \in \boldsymbol{U}} \leftarrow$ random initialisation
3: **while** $P(\mathcal{D}|\{\theta_U^{t+1}\}_{U \in \boldsymbol{U}}) \geq P(\mathcal{D}|\{\theta_U^t\}_{U \in \boldsymbol{U}})$ **do**
4:     **for** $U \in \boldsymbol{U}$ **do**
5:         $\theta_U^{t+1} \leftarrow |\mathcal{D}|^{-1} \sum_{\boldsymbol{x} \in \mathcal{D}} \theta_{U|\boldsymbol{x}}^t$
6:         $t \leftarrow t + 1$
7:     **end for**
8: **end while**
9: **return** $\{\theta_U^{t+1}\}_{U \in \boldsymbol{U}}$

FSCM (=BN) QUERIES

This is a single run, returning exogenous chances
to be iterated for different random initialisations

# Speeding up the Causal EM

- Parallelisation (on multiple levels)

  - EM initialisations

  - Dataset records

  - (Connected Components)

- Knowledge Compilation?

- EM queries on different models

  - initialisation $\theta_0$

  - iteration $t$

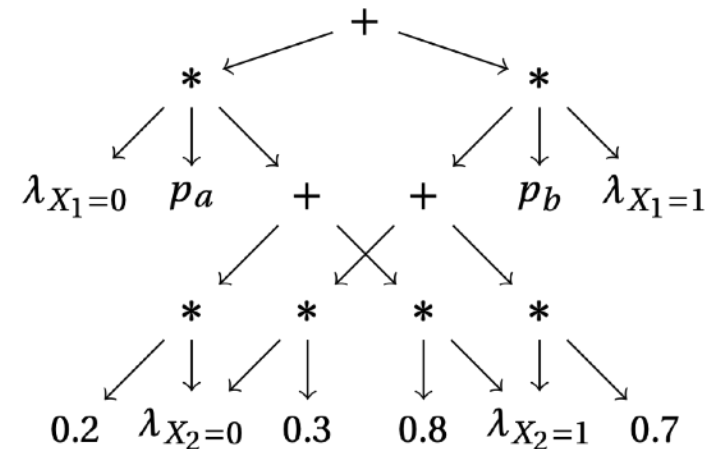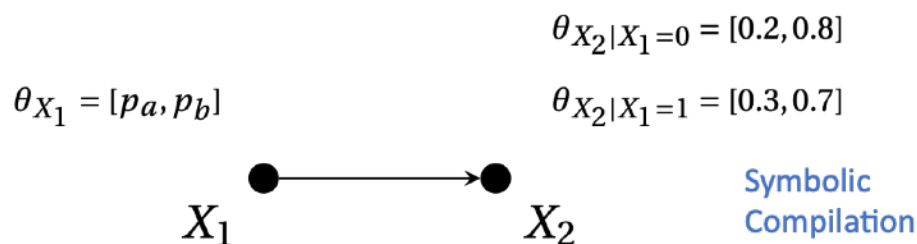- Multiple compilations could be expensive, but ...

1: $t \leftarrow 0$
2: **while** $P(\mathcal{D}|\{\theta_U^{t+1}\}_{U \in U}) \geq P(\mathcal{D}|\{\theta_U^t\}_{U \in U})$ **do**
3:    **for** $U \in U$ **do**
4:       $\theta_U^{t+1} \leftarrow |\mathcal{D}|^{-1} \sum_{v \in \mathcal{D}} \theta_{U|v}^t$
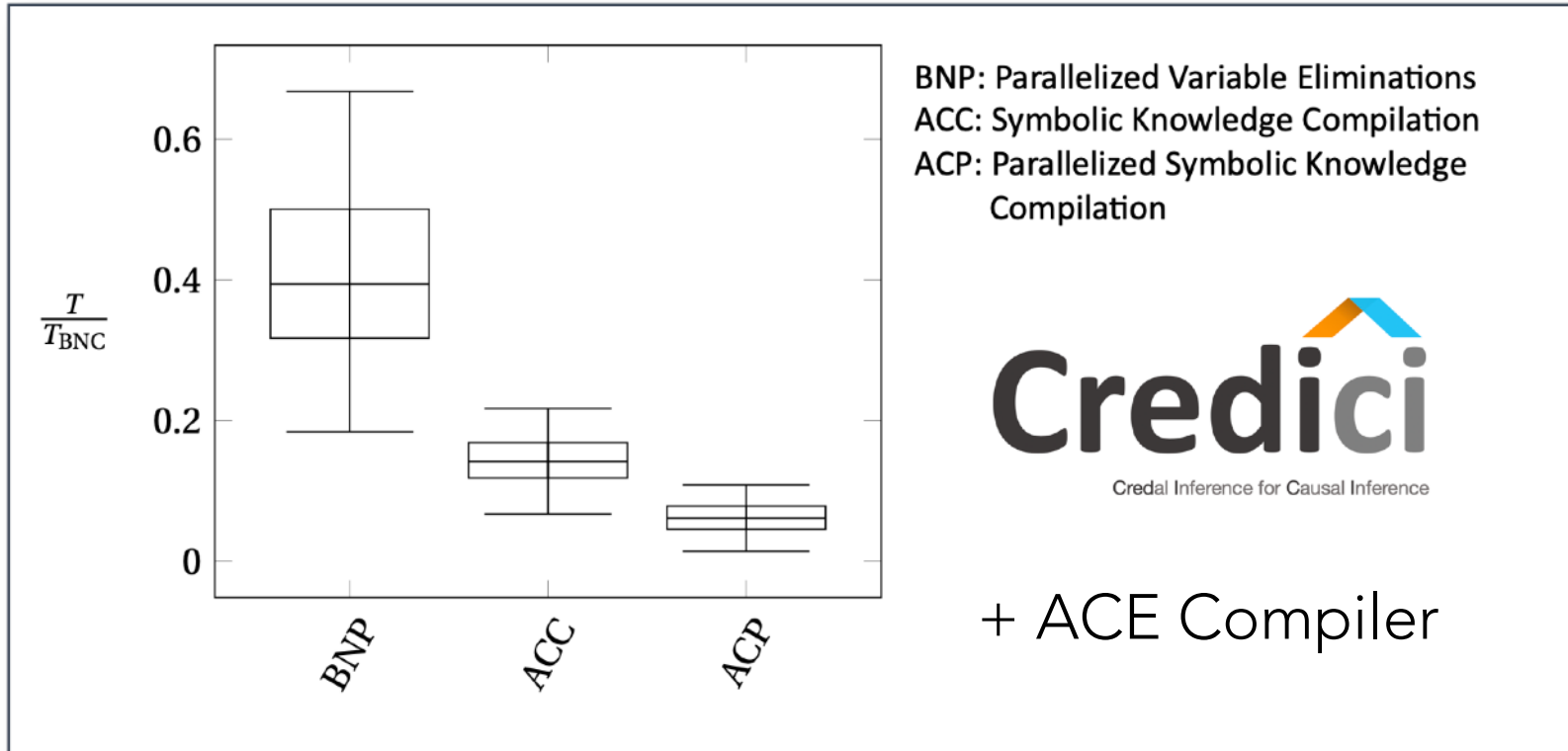5:       $t \leftarrow t+1$
6:    **end for**
7: **end while**

# Symbolic Knowledge Compilation

- Multiple inferences on different FSCM models
- All FSCMs have a shared structure:
  - Same variables and graph
  - Same equations (endogenous CPTs)
- A "symbolic" (parametrised) compilation
- A single compilation with unique parameters (used as IDs)
- Re-compilation by changing the parameters (linear time wrt pars)



$\theta_{X_2|X_1=0} = [0.2, 0.8]$

$\theta_{X_1} = [p_a, p_b]$     $\theta_{X_2|X_1=1} = [0.3, 0.7]$

$X_1 \longrightarrow X_2$     **Symbolic Compilation**

# Preliminary Experiments



BNP: Parallelized Variable Eliminations
ACC: Symbolic Knowledge Compilation
ACP: Parallelized Symbolic Knowledge
       Compilation

**Credici**
Credal Inference for Causal Inference

+ ACE Compiler

- Symbolic compilation more effective than (component) parallelisation
- ACE exploits the determinism in the structural equations
- Overall, one order of magnitude faster with parallelisation + KC

# Conclusions and (a Lot of) Future Work

- Conclusions
    - Concept of parametrised compilation of circuits
    - Knowledge compilation to tractable arithmetic circuits achieves SOTA performance in counterfactual bounding
- Future Work
    - Specialised compilation for SCMs? Canonical equations (FO?), connected components (Decomposed?) and counterfactual graphs (Lifted Inference?)
    - Query-aware methods? (current are query-agnostic)
    - Genuine symbolic inference ("credal" causal EM)
    - Better parallelisation (Julia)