

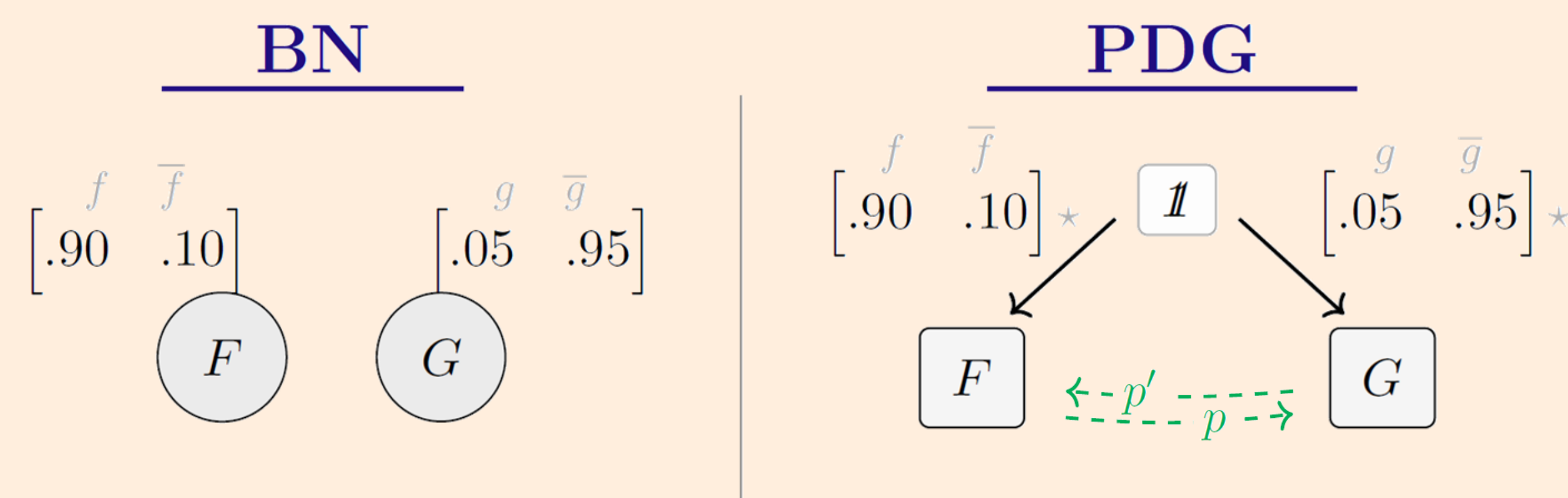
WHY YET ANOTHER GRAPHICAL MODEL?

PDGs...

- capture inconsistency, including conflicting information from multiple sources with varying reliability.
- are especially modular; to combine info from two sources, simply take a PDG union. This incorporates new data (edge cpds) and concepts (nodes) without affecting previous information.
- cleanly separate quantitative info (the cpds) from qualitative info (the edges), with variable confidence in both (the weights β and α). This is captured by terms *Inc* and *IDef* in our scoring function.
- have (several) natural semantics; one of them allows us to pick out a unique distribution. Using this distribution, PDGs can capture BNs and factor graphs.

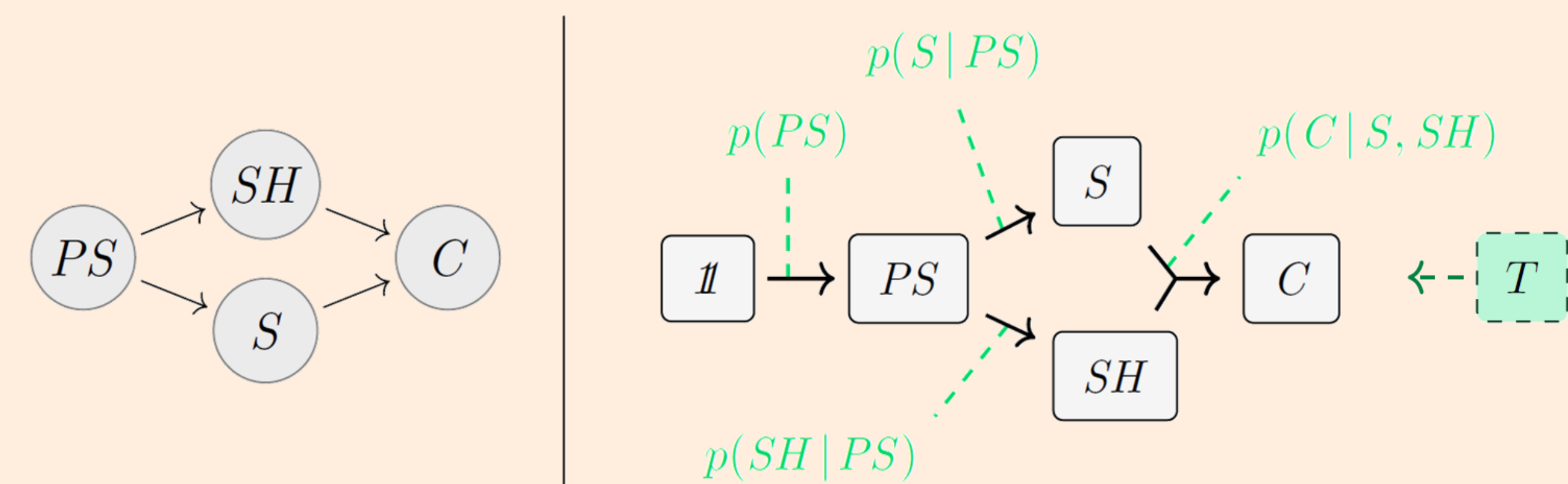
MODELING EXAMPLES

A SIMPLE ILLUSTRATION



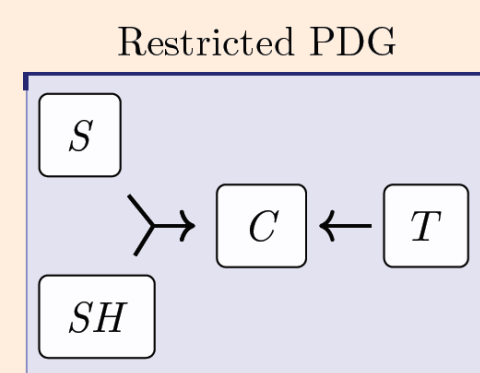
- The cpds of a PDG are attached to edges, not nodes.
- PDGs can incorporate arbitrary new probabilistic information.
- PDGs can be inconsistent
 - ...but BNs must resolve inconsistency first, which may break symmetry and irreversibly lose information.

BAYESIAN NETWORKS AS PDGs

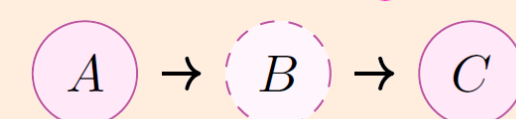


In contrast with BNs:

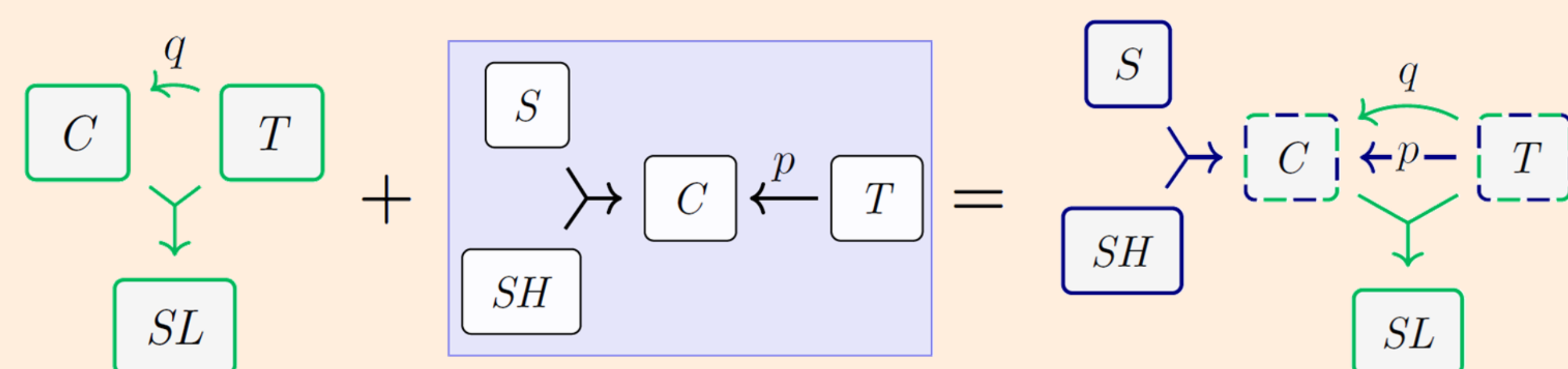
- edge composition has *quantitative* meaning, since edges have cpds;
- a variable can be the target of more than one cpd;
- arbitrary restrictions of PDGs are still PDGs.



In a qualitative BN: removing data results in *new* knowledge: $A \perp C$.



COMBINING PDGs



- Arbitrary PDGs may be combined without loss of information
- They may have parallel edges (e.g., p, q), which directly conflict.

PROBABILISTIC DEPENDENCY GRAPHS

Oliver E. Richardson Joseph Y. Halpern

Cornell University
Department of Computer Science

Definition (Probabilistic Dependency Graph)

A PDG is a tuple $\mathbf{m} = (\mathcal{N}, \mathcal{E}, \mathcal{V}, \mathbf{p}, \alpha, \beta)$, where

\mathcal{N} is a finite set of nodes (variables)

\mathcal{V} gives a set $\mathcal{V}(X)$ of possible values for each X ;

\mathcal{E} is a set of labeled edges $\{X \xrightarrow{L} Y\}$,

and associated to each $X \xrightarrow{L} Y$, there is:

\mathbf{p}_L a cpd $\mathbf{p}_L(Y | X)$;

$\alpha_L \in [0, \infty)$ a confidence in the functional dependence $X \rightarrow Y$

$\beta_L \in (0, \infty)$ a confidence in the reliability of \mathbf{p}_L .

PDG SEMANTICS

- $\{\mathbf{m}\}$ The set of joint distributions consistent with \mathbf{m} ;
- $\llbracket \mathbf{m} \rrbracket_\gamma$ A loss function (parameterized by γ), scoring a joint distribution's compatibility with \mathbf{m} ;

$$\llbracket \mathbf{m} \rrbracket_\gamma(\mu) := \text{Inc}_m(\mu) + \gamma \text{IDef}_m(\mu)$$

tradeoff parameter $\gamma \geq 0$

(Quantitative)

Definition (Inc)

The *incompatibility* of μ with \mathbf{m} :

$$\text{Inc}_m(\mu) := \sum_{X \xrightarrow{L} Y} \beta_L \mathbf{D}(\mu_{Y|X} \| \mathbf{p}_L)$$

The *inconsistency* of \mathbf{m} is

$$\text{Inc}(\mathbf{m}) := \inf_{\mu \in \Delta \mathcal{V}(\mathbf{m})} \text{Inc}_m(\mu).$$

(Qualitative)

Definition (IDef)

The *m-information deficit* of μ :

bits to separately determine each target, knowing the source

$$\text{IDef}_m(\mu) = \sum_{X \xrightarrow{L} Y} \alpha_L \mathbf{H}_\mu(Y | X) - \mathbf{H}(\mu)$$

bits to determine all vars

- $\llbracket \mathbf{m} \rrbracket^*$ The (unique) "best" joint distribution (in the quantitative limit).

$$\llbracket \mathbf{m} \rrbracket^* := \lim_{\gamma \rightarrow 0} \arg \min_{\mu} \llbracket \mathbf{m} \rrbracket_\gamma(\mu)$$

PROPERTIES OF SEMANTICS

Proposition (the second semantics extends the first)

$$\llbracket \mathbf{m} \rrbracket = \{\mu : \llbracket \mathbf{m} \rrbracket_0(\mu) = 0\}.$$

Proposition (If there are distributions consistent with \mathbf{m} , the best distribution is one of them.)

$$\llbracket \mathbf{m} \rrbracket^* \in \llbracket \mathbf{m} \rrbracket_0, \text{ so if } \mathbf{m} \text{ is consistent, then } \llbracket \mathbf{m} \rrbracket^* \in \llbracket \mathbf{m} \rrbracket.$$

Proposition (uniqueness for small γ)

- If $0 < \gamma \leq \min_L \beta_L^m$, then $\llbracket \mathbf{m} \rrbracket_\gamma^*$ is a singleton.
- $\lim_{\gamma \rightarrow 0} \llbracket \mathbf{m} \rrbracket_\gamma^*$ exists and is unique.

CAPTURING BNs AS PDGs

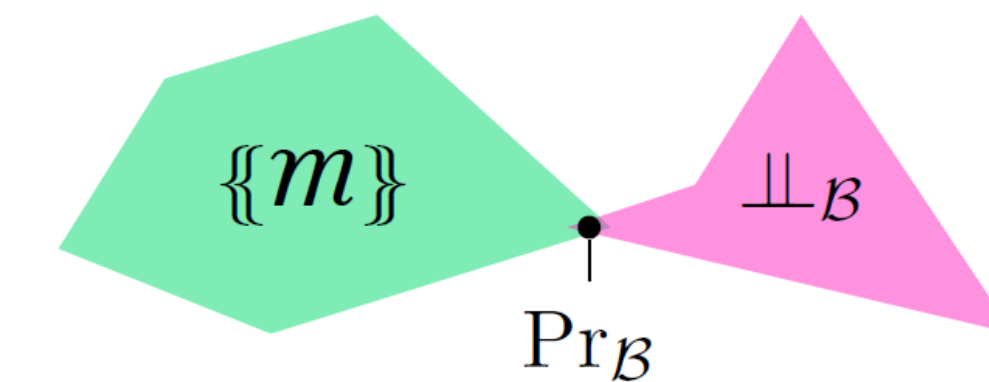
Let $\mathbf{m}_{\mathcal{B}, \beta}$ be the PDG corresponding to the BN \mathcal{B} , with weights β .

Theorem (BNs are PDGs)

If \mathcal{B} is a BN and $\text{Pr}_{\mathcal{B}}$ is the distribution it specifies, then for all $\gamma > 0$ and all vectors β ,

$$\llbracket \mathbf{m}_{\mathcal{B}, \beta} \rrbracket_\gamma^* = \text{Pr}_{\mathcal{B}}.$$

space of distributions consistent with $\mathbf{m}_{\mathcal{B}}$ (which minimize *Inc*)

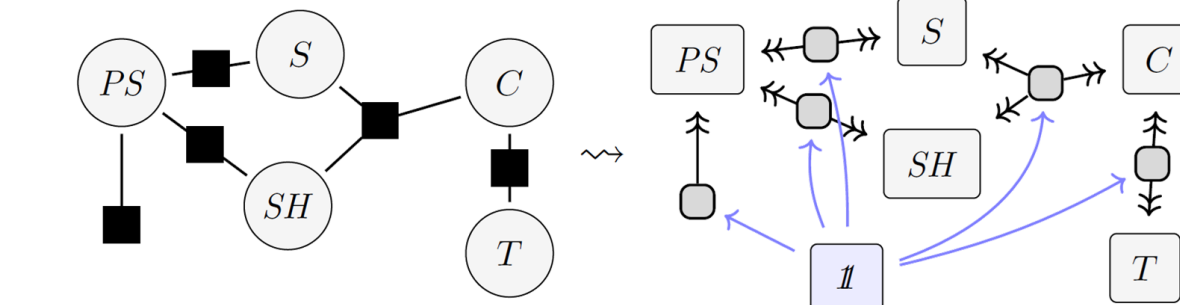


space of distributions with independencies of \mathcal{B} (which can be shown to minimize *IDef*)

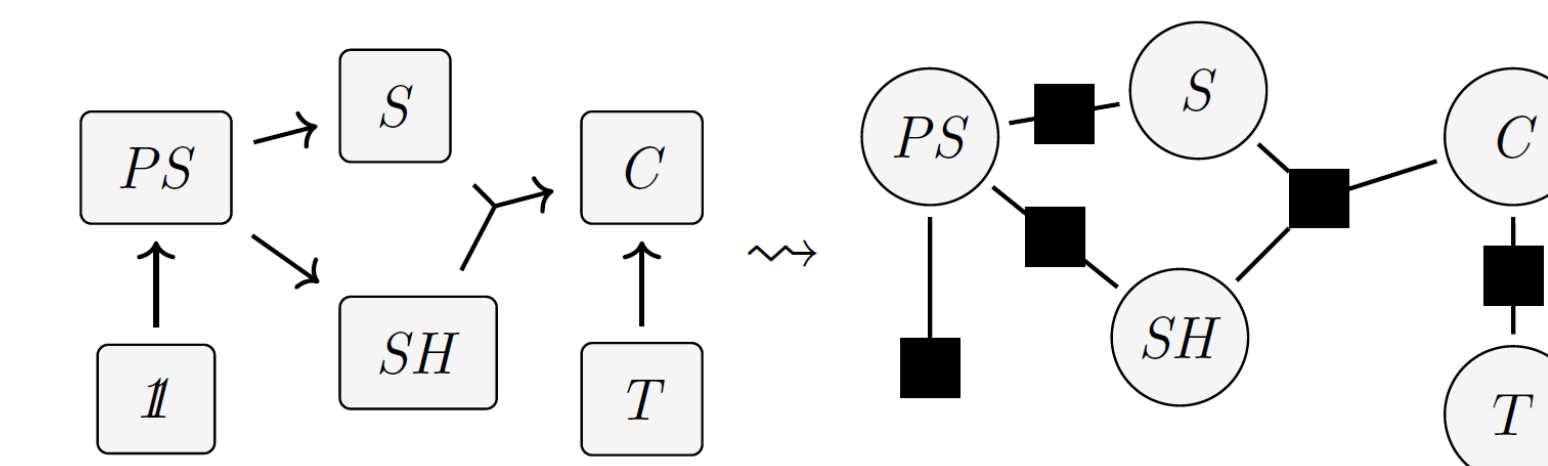
CAPTURING FACTOR GRAPHS AS PDGs

Theorem (PDGs capture factor graphs)

We can naturally translate factor graphs and their exponential families, into PDGs, in a way which preserves their semantics.



PDGs AS FACTOR GRAPHS



The cpds of a PDG are essentially factors. Are the semantics the same? Only for $\gamma = 1$.

Theorem

$\llbracket \mathbf{n} \rrbracket_1^* = \text{Pr}_{\Phi_n}$ for all unweighted PDGs \mathbf{n} .

$$\llbracket \mathbf{m} \rrbracket_\gamma(\mu) = \mathbb{E}_{\mathbf{w} \sim \mu} \left\{ \sum_{X \xrightarrow{L} Y} \left[\beta_L \log \frac{1}{\mathbf{p}_L(y^{\mathbf{w}} | x^{\mathbf{w}})} + (\alpha_L \gamma - \beta_L) \log \frac{1}{\mu(y^{\mathbf{w}} | x^{\mathbf{w}})} \right] - \gamma \log \frac{1}{\mu(\mathbf{w})} \right\}.$$

local regularization ($\beta_L > \alpha_L \gamma$) global regularization

INFERENCE AND INCONSISTENCY: A GLIMPSE.

Conditioning as inconsistency resolution.

To condition on $Y = y$, in \mathbf{m} , simply add the edge $\mathbb{I} \xrightarrow{\delta_y} Y$ to get $\mathbf{m}_{Y=y}$. Then $\llbracket \mathbf{m}_{Y=y} \rrbracket^* = \llbracket \mathbf{m} \rrbracket^* | (Y = y)$.

Querying $\text{Pr}(Y | X)$ in a PDG \mathbf{m} .

- We can add $X \xrightarrow{p} Y$ to \mathbf{m} with a cpt p , to get \mathbf{m}^{+p} .
- The choice of cpd p that minimizes the inconsistency of \mathbf{m}^{+p} (which is strongly convex and smooth in p) is $\llbracket \mathbf{m} \rrbracket^*(Y | X)$,
- so oracle access to inconsistency yields fast inference by gradient descent.

