

***StarfishDB:***  
*Probabilistic Programming*  
*Datalog in Action*

Niccolò Meneghetti  
*University of Michigan-Dearborn*  
niccolom@umich.edu

# Outline

## **StarfishDB: a Query Execution Engine for Relational Probabilistic Programming**

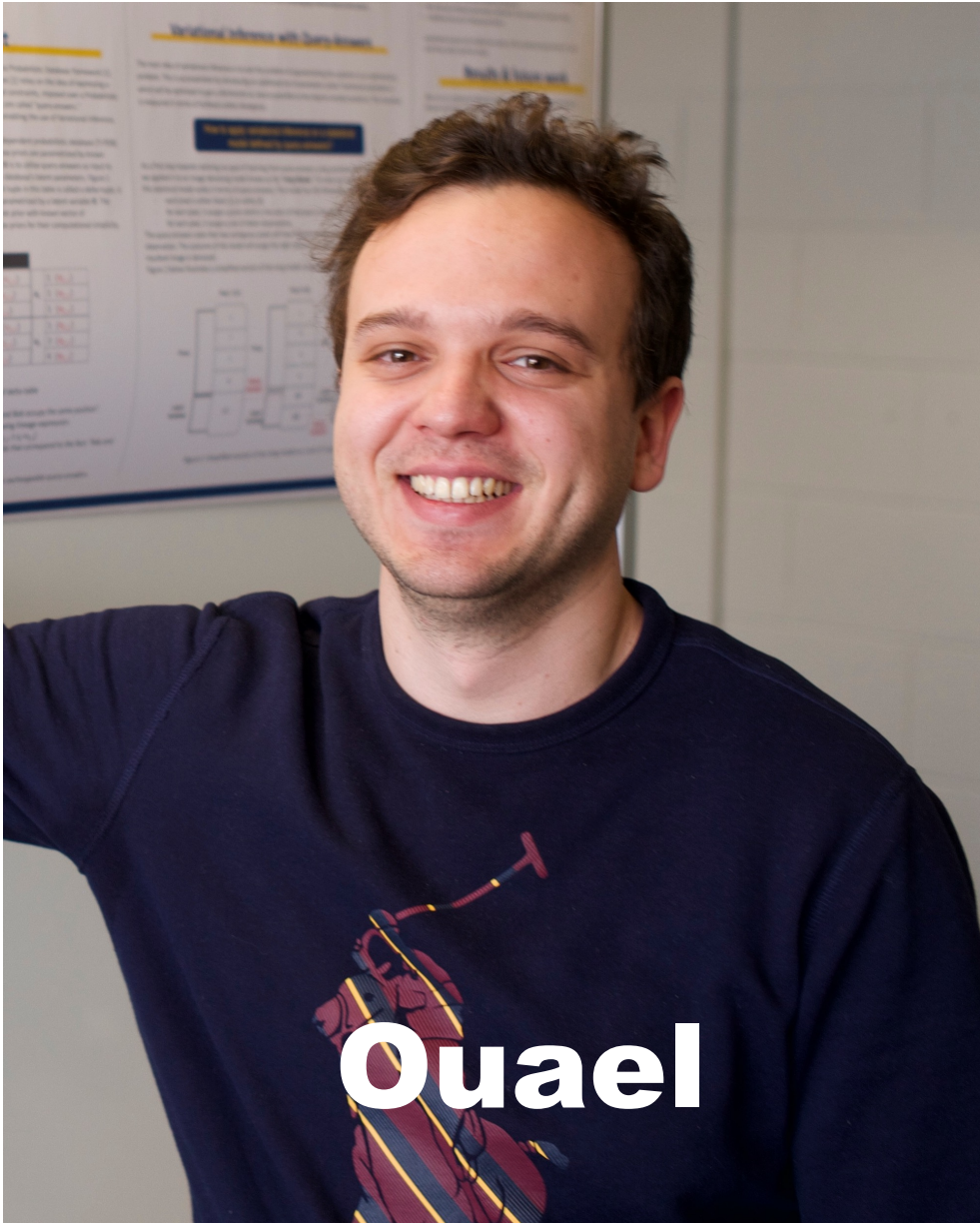
OUAEL BEN AMARA\*, University of Michigan-Dearborn, U.S.A.

SAMI HADOUAJ\*, University of Michigan-Dearborn, U.S.A.

NICCOLÒ MENEGHETTI, University of Michigan-Dearborn, U.S.A.

- De Finetti Logic Programming
- Variational Inference
- Probabilistic Programming Datalog

# The Authors

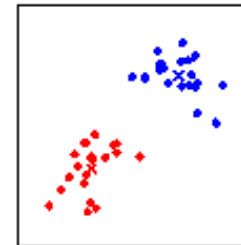
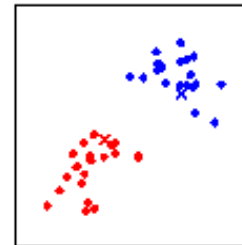
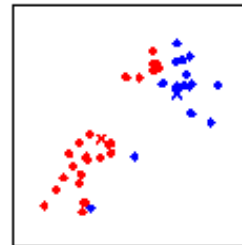
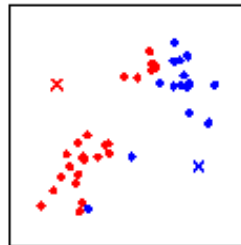
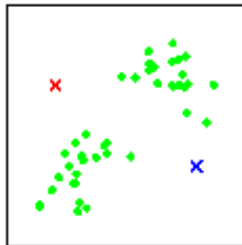
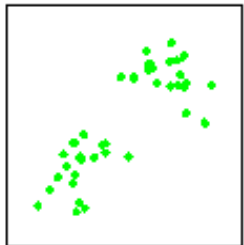


# Probabilistic Programming

**Generative Story:**  $z \sim \text{Categorical}(\phi)$   
 $\mathbf{x} \sim \text{Gaussian}(\mu_z, \Sigma_z)$

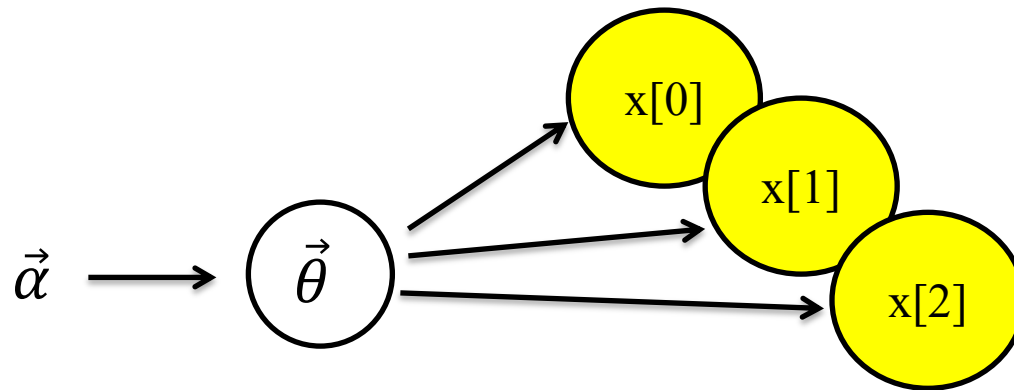
**Data:**  $\mathcal{D} = \{\mathbf{x}^{(i)}\}_{i=1}^N$  where  $\mathbf{x}^{(i)} \in \mathbb{R}^M$

Goal: compute the **posterior density** of the generative process w.r.t. the data.



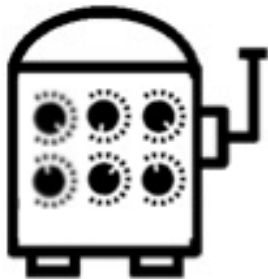
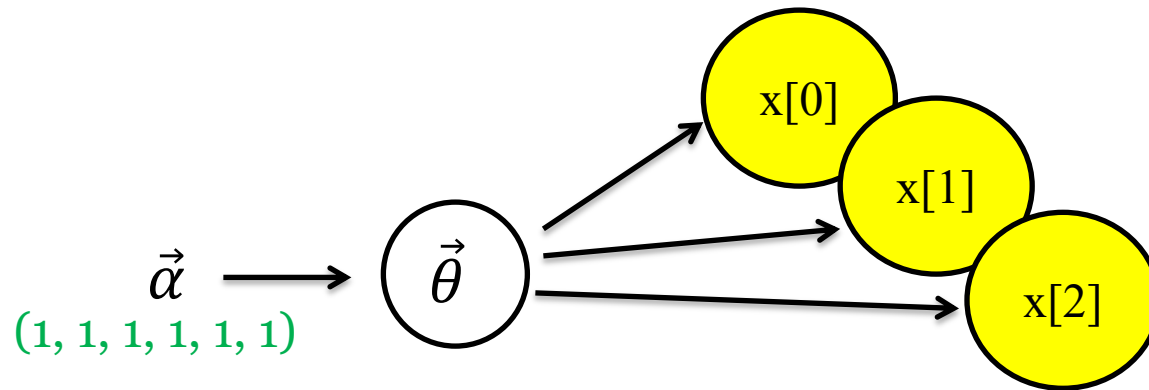
# Pólya Dice

**Pólya die** → a *Categorical* distribution (parametrized by  $\vec{\theta}$ )  
with a *Dirichlet* prior (parametrized by  $\vec{\alpha}$ )



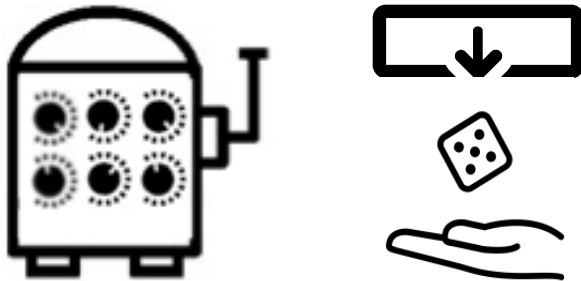
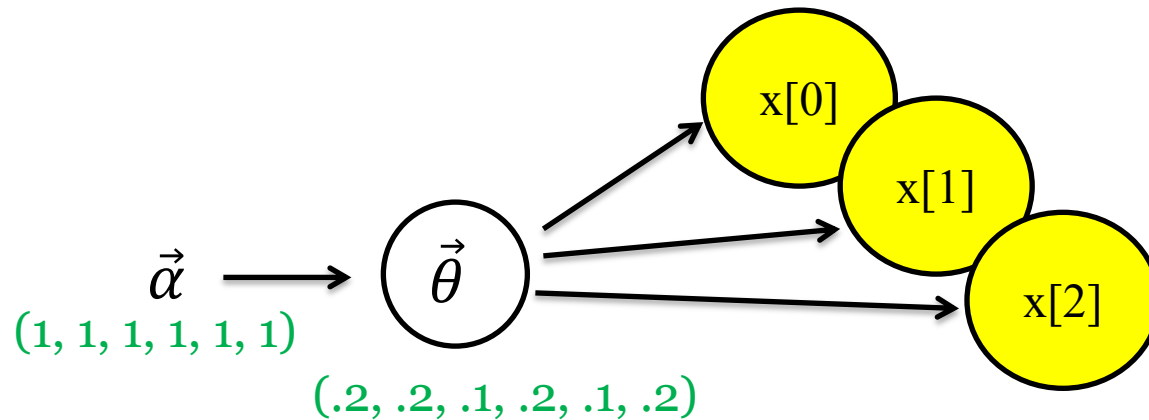
# Pólya Dice

**Pólya die** → a *Categorical* distribution (parametrized by  $\vec{\theta}$ )  
with a *Dirichlet* prior (parametrized by  $\vec{\alpha}$ )



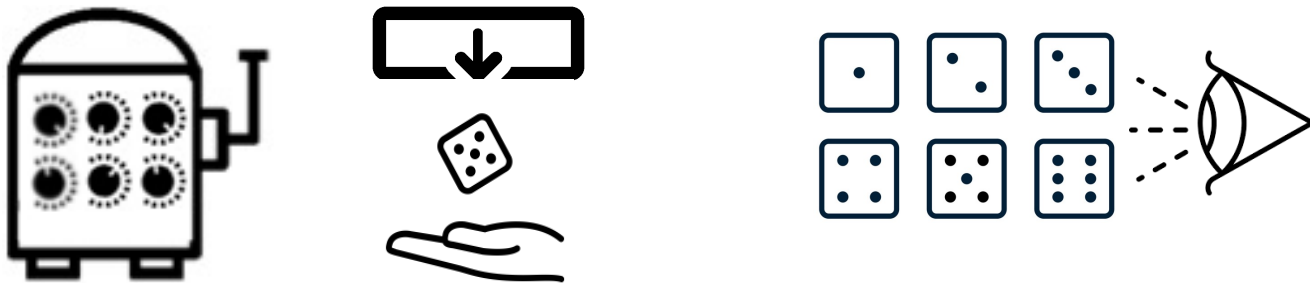
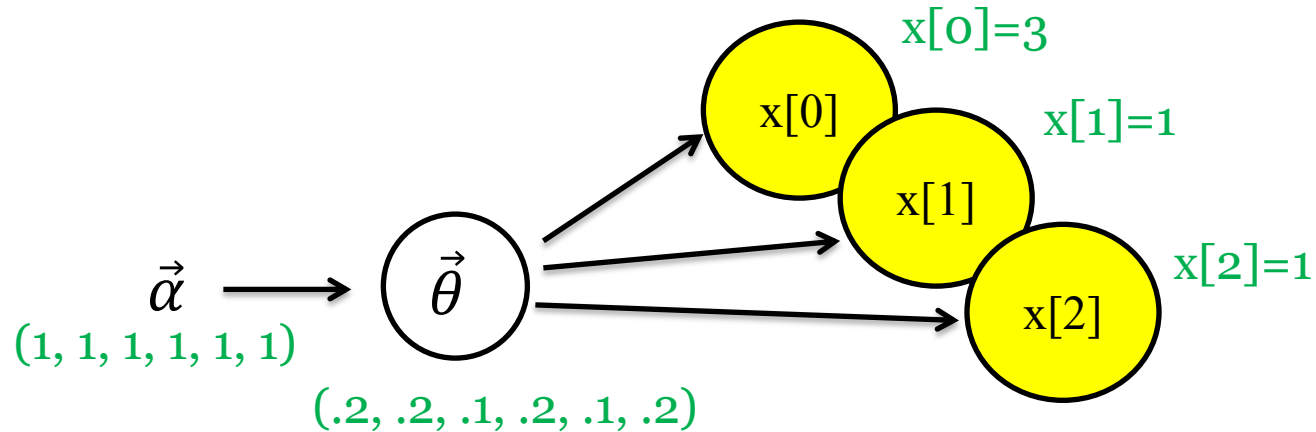
# Pólya Dice

Pólya die  $\rightarrow$  a *Categorical* distribution (parametrized by  $\vec{\theta}$ )  
with a *Dirichlet* prior (parametrized by  $\vec{\alpha}$ )



# Pólya Dice

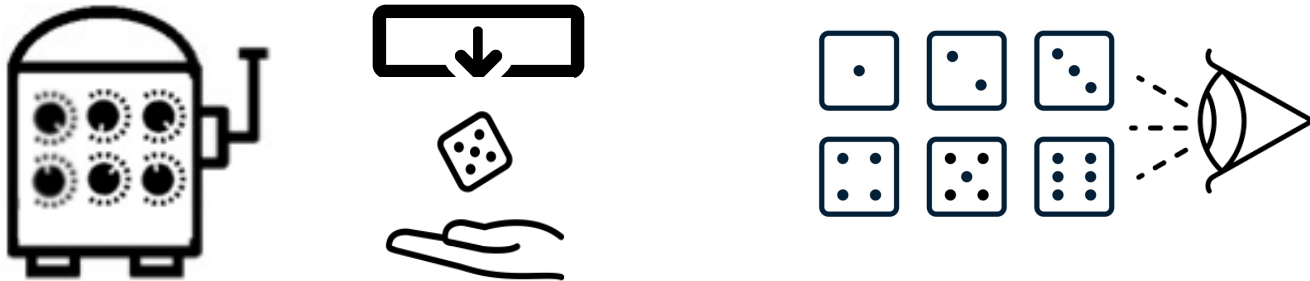
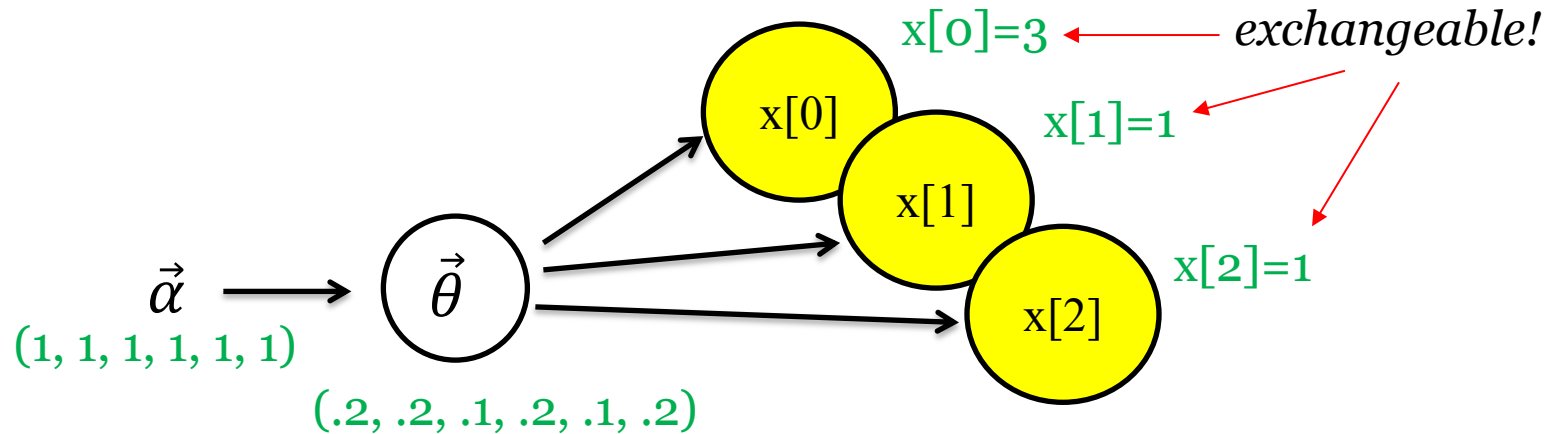
**Pólya die** → a *Categorical* distribution (parametrized by  $\vec{\theta}$ )  
with a *Dirichlet* prior (parametrized by  $\vec{\alpha}$ )





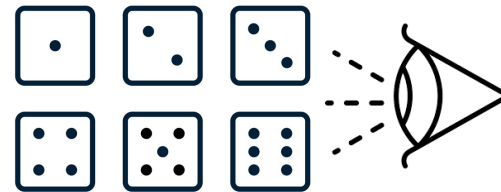
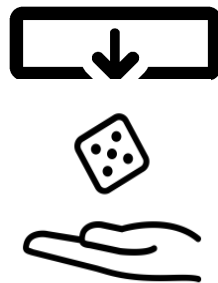
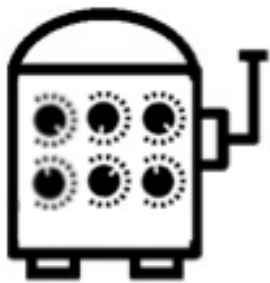
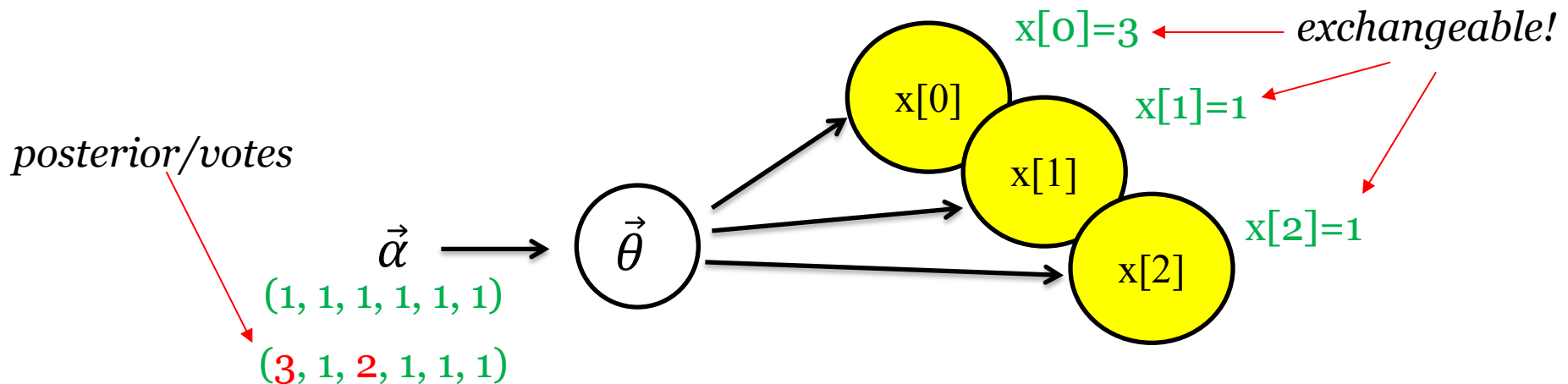
# Pólya Dice

**Pólya die** → a *Categorical* distribution (parametrized by  $\vec{\theta}$ )  
with a *Dirichlet* prior (parametrized by  $\vec{\alpha}$ )



# Pólya Dice

**Pólya die** → a *Categorical* distribution (parametrized by  $\vec{\theta}$ )  
with a *Dirichlet* prior (parametrized by  $\vec{\alpha}$ )

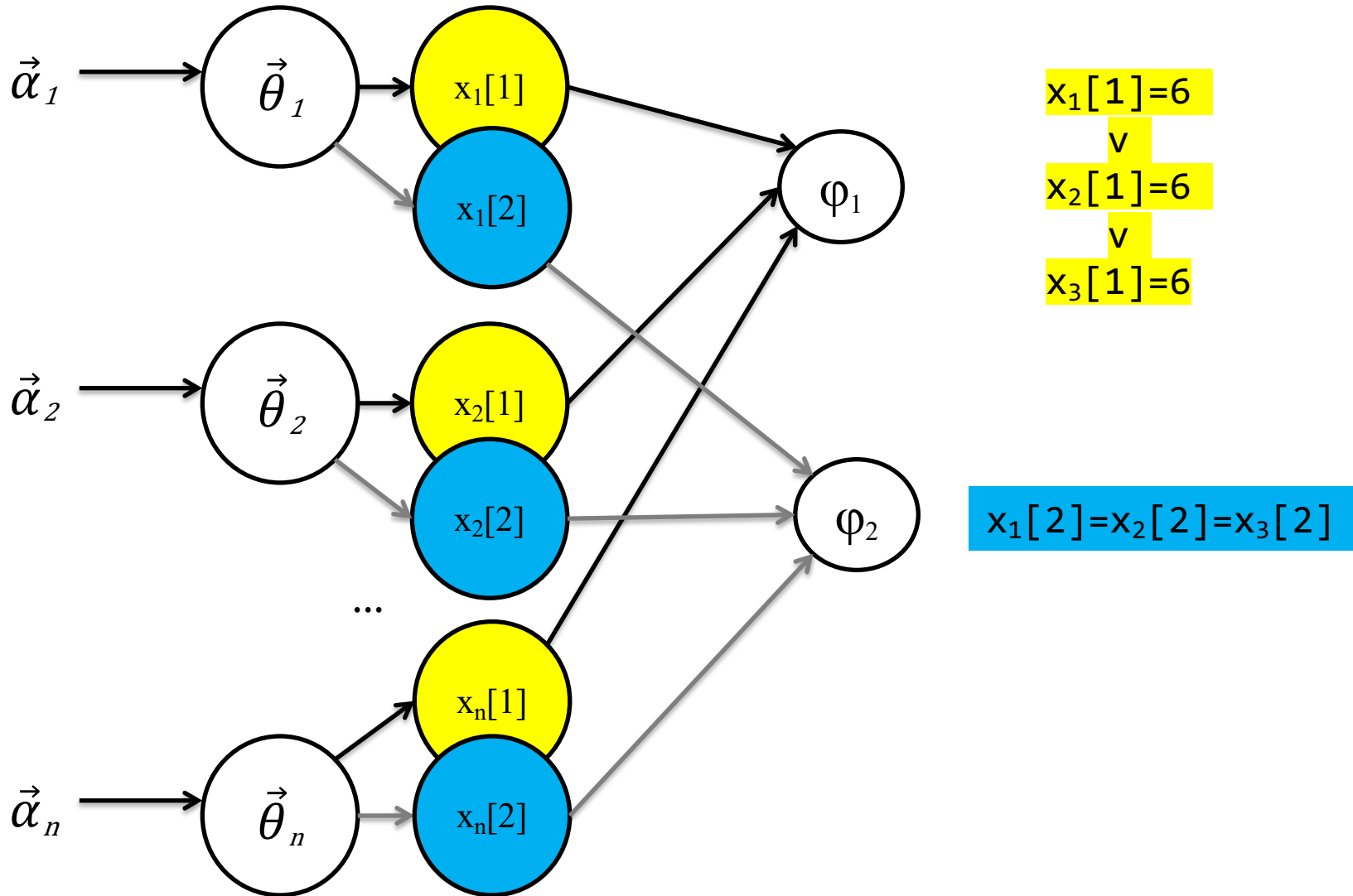


# De Finetti Logic Programming

A *DFL Program* consists of:

1. A generative process  
(defined as a set of pairwise independent Pólya dice)
2. A set of exchangeable constraints

# Exchangeable Constraints



# De Finetti Logic Programming

A *DFL Program* consists of:

1. A generative process  
(defined as a set of pairwise independent Pólya dice)
2. A set of exchangeable constraints

$$p(\Theta \mid \Phi, \mathbb{A}) = \sum_{\tau \in \text{SAT}(\Phi, \mathbb{X})} p(\Theta \mid \tau, \mathbb{A}) \cdot P(\tau \mid \Phi, \mathbb{A})$$

*product of Dirichlet densities*

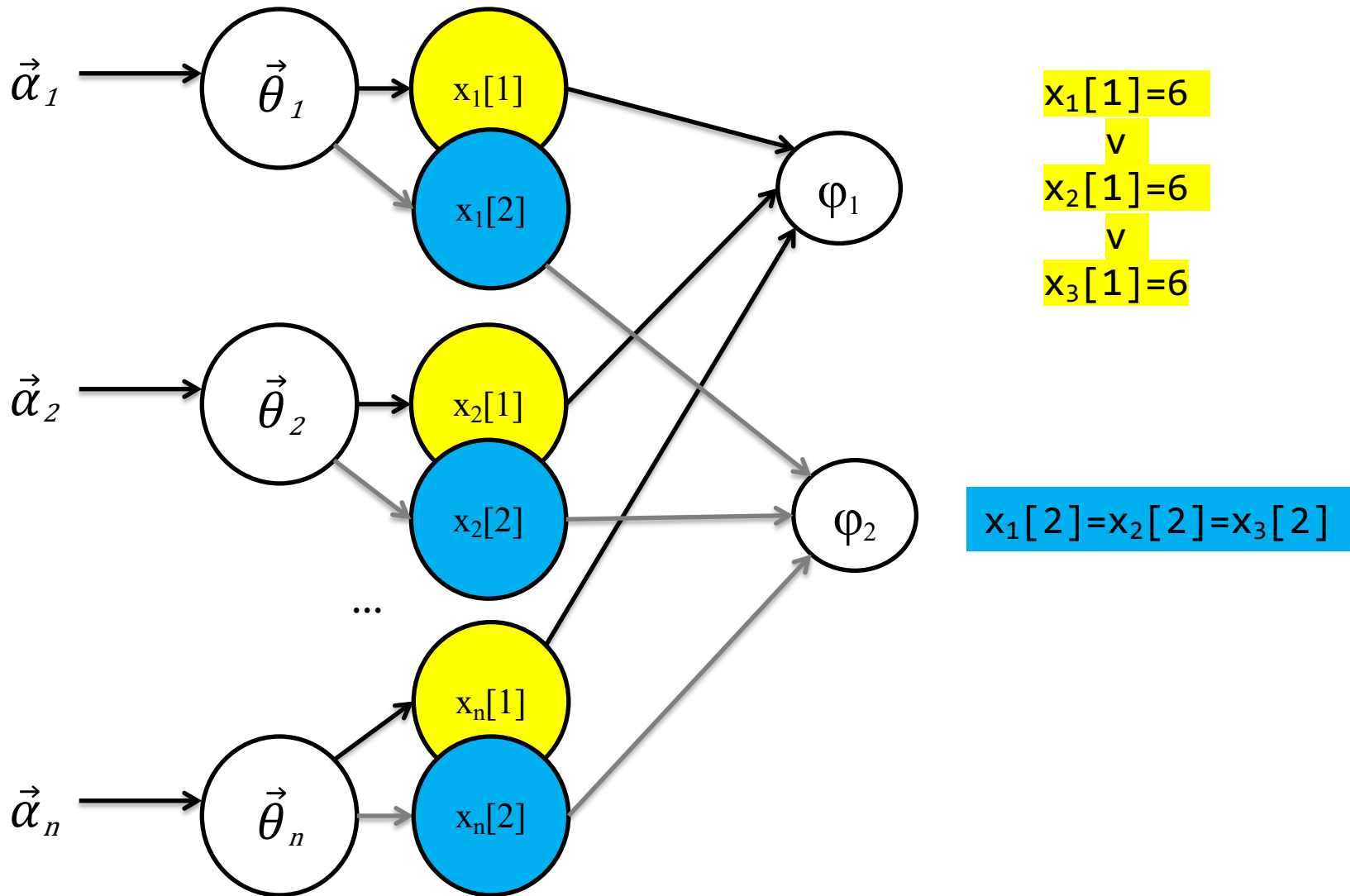
posterior

assignments to  $X$   
that satisfy  $\Phi$

Pólya urn

The diagram illustrates the decomposition of a posterior distribution. The equation shows the posterior  $p(\Theta \mid \Phi, \mathbb{A})$  as a sum over all assignments  $\tau$  to  $X$  that satisfy the constraints  $\Phi$ . Each assignment  $\tau$  contributes a term  $p(\Theta \mid \tau, \mathbb{A}) \cdot P(\tau \mid \Phi, \mathbb{A})$ . The term  $p(\Theta \mid \tau, \mathbb{A})$  is identified as the 'product of Dirichlet densities', and  $P(\tau \mid \Phi, \mathbb{A})$  is identified as the 'Pólya urn' distribution. Red arrows point from the labels to the corresponding parts of the equation.

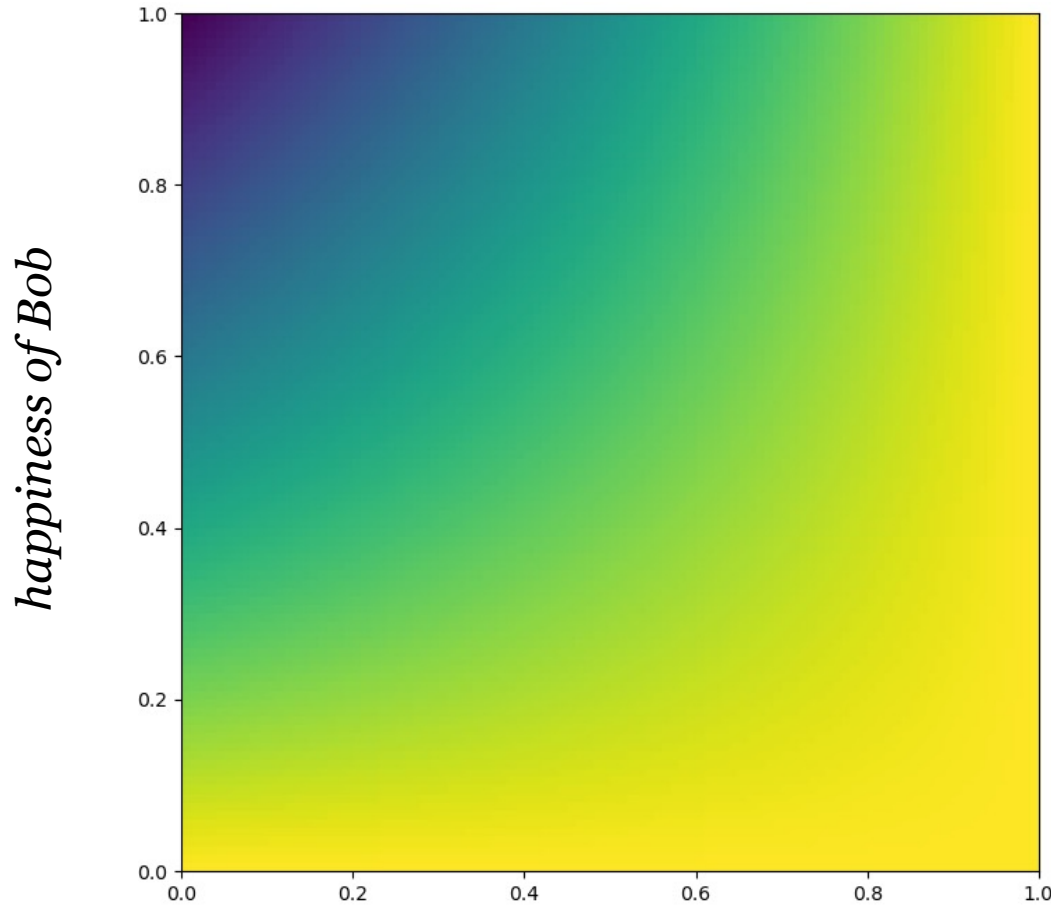
# Exchangeable Constraints



# Posterior of two Pólya coins

Ada ☹️  
Bob 😊

Ada 😊  
Bob 😊



**C1:** If Bob is happy, so is Ada

Ada ☹️  
Bob 😊

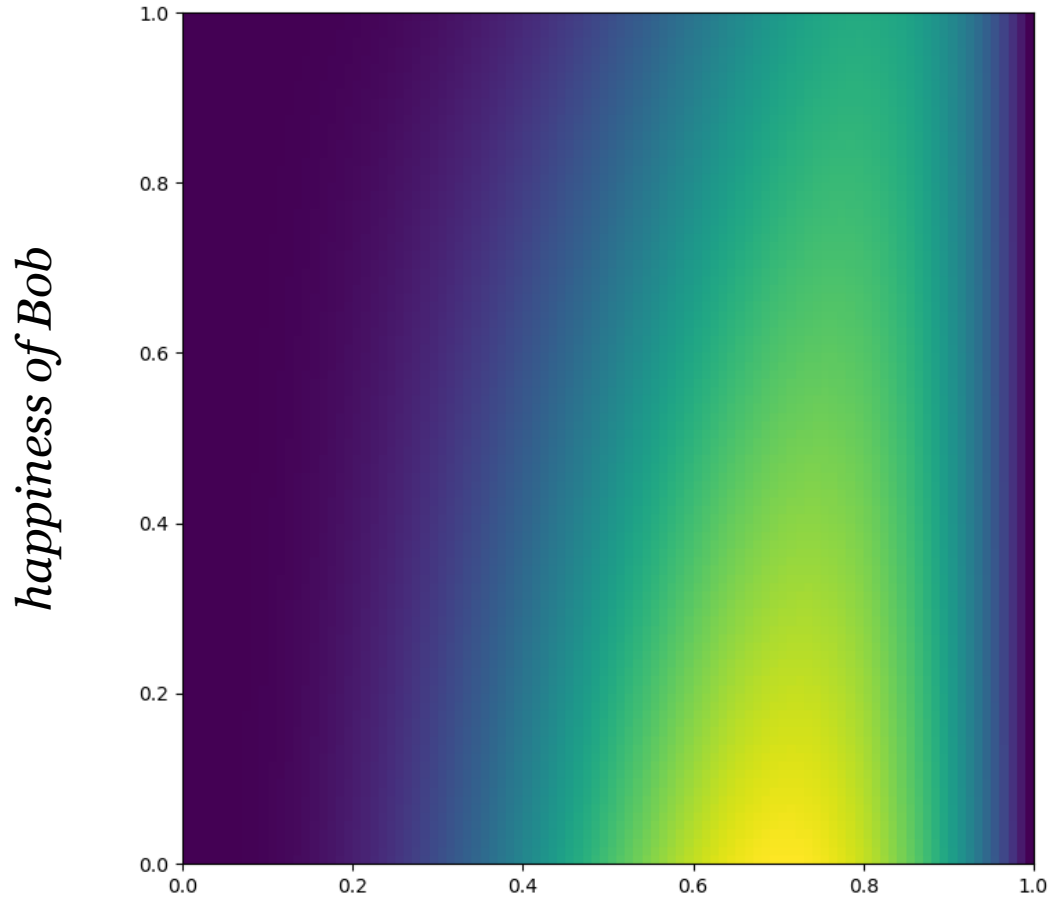
*happiness of Ada*

Ada 😊  
Bob ☹️

# Posterior of two Pólya coins

Ada 😞  
Bob 😊

Ada 😊  
Bob 😊



- C1: If Bob is happy, so is Ada
- C2: Ada is happy
- C3: Ada is happy
- C4: Ada is happy
- C5: Ada is sad

Ada 😞  
Bob 😊

happiness of Ada

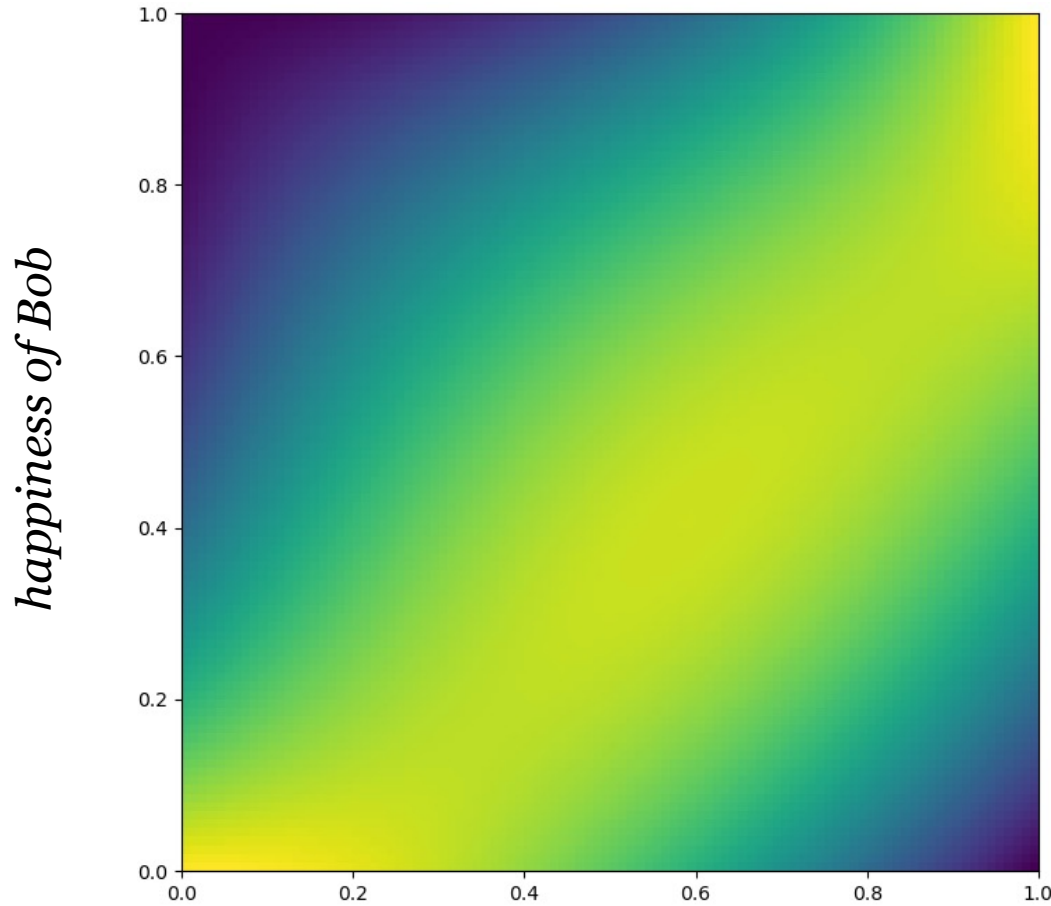
Ada 😊  
Bob 😞



# Posterior of two Pólya coins

Ada 😞  
Bob 😊

Ada 😊  
Bob 😊



**C1:** If Bob is happy, so is Ada  
**C2:** Ada and Bob have the same mood

Ada 😞  
Bob 😊

*happiness of Ada*

Ada 😊  
Bob 😞

# De Finetti Logic Programming

Computing the exact posterior is impractical. Two solutions:

1. Gibbs Sampling
2. Variational Inference

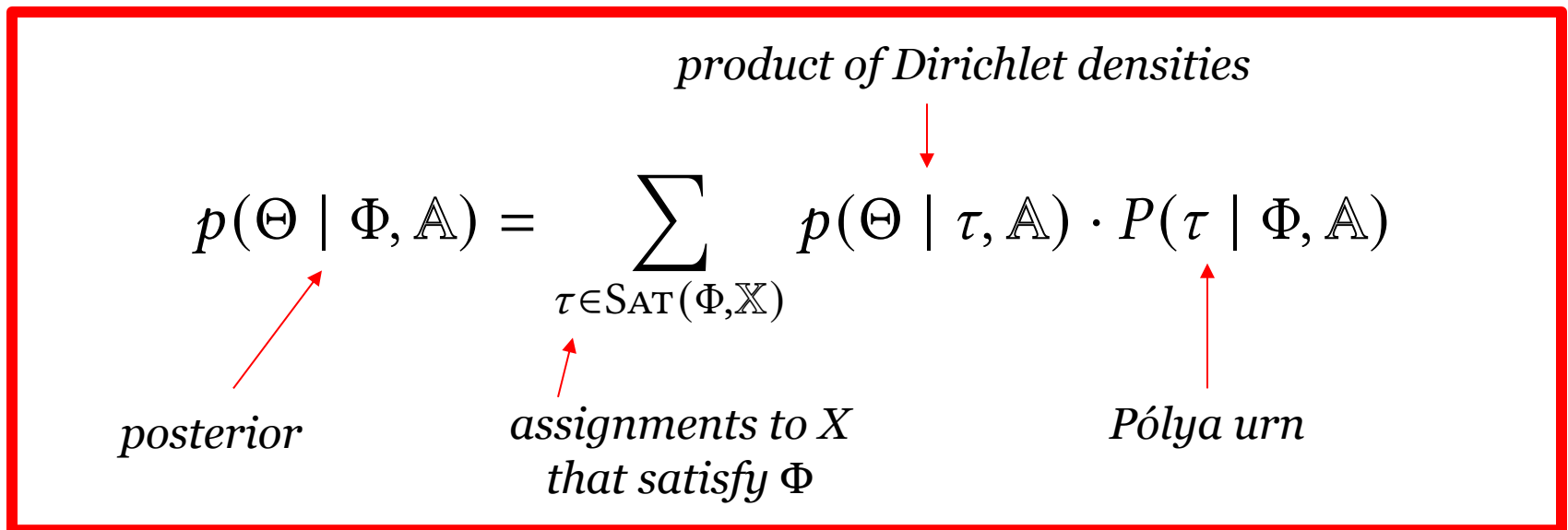
$$p(\Theta \mid \Phi, \mathbb{A}) = \sum_{\tau \in \text{SAT}(\Phi, \mathbb{X})} p(\Theta \mid \tau, \mathbb{A}) \cdot P(\tau \mid \Phi, \mathbb{A})$$

*product of Dirichlet densities*

*posterior*

*assignments to  $X$   
that satisfy  $\Phi$*

*Pólya urn*



# Variational Inference

*variational approximation*  $q(\Theta, \Phi) \stackrel{\text{def}}{=} q_{\Theta} \cdot q_{\Phi}$

$q_{\Theta} \stackrel{\text{def}}{=} \prod_{\theta_i \in \Theta} q_{\theta_i}(\theta_i | \mu_i)$  *Dirichlet densities*

$q_{\Phi} \stackrel{\text{def}}{=} \prod_{\phi_m \in \Phi} q_{\phi_m}(\tau | \lambda_m)$  *Categorical distributions*

*product of Dirichlet densities*

*posterior*  $p(\Theta | \Phi, \mathbb{A}) = \sum_{\tau \in \text{SAT}(\Phi, \mathbb{X})} p(\Theta | \tau, \mathbb{A}) \cdot P(\tau | \Phi, \mathbb{A})$

*assignments to  $X$  that satisfy  $\Phi$*

*Pólya urn*

# Variational Inference

*variational approximation*  $\nearrow$

$$q(\Theta, \Phi) \stackrel{\text{def}}{=} q_{\Theta} \cdot q_{\Phi}$$
$$q_{\Theta} \stackrel{\text{def}}{=} \prod_{\theta_i \in \Theta} q_{\theta_i}(\theta_i \mid \boldsymbol{\mu}_i) \quad \leftarrow \text{Dirichlet densities}$$
$$q_{\Phi} \stackrel{\text{def}}{=} \prod_{\phi_m \in \Phi} q_{\phi_m}(\tau \mid \boldsymbol{\lambda}_m) \quad \leftarrow \text{Categorical distributions}$$

*optimize  $\boldsymbol{\mu}_i$*   $\longrightarrow$   $q_{\theta_i}(\theta_{i,v}) \propto \boldsymbol{\alpha}_i + \sum_{\phi_m \in \Phi} \left[ \sum_{\tau \in \text{SAT}(\phi_m \wedge (X_i=v))} q_{\phi_m}(\tau) \right]$

*optimize  $\boldsymbol{\lambda}_m$*   $\longrightarrow$   $q_{\phi_m}(\tau) \propto \prod_{\theta_i \in \Theta} \exp \mathbb{E}_q \left[ \log q_{\theta_i}(\theta_{i,\tau[i]}) \right]$

# Where do the constraints come from?

EMP	ROLE
Ada ( $x_1$ )	Lead ( $v_{1,1}$ )
	Dev ( $v_{1,2}$ )
	QA ( $v_{1,3}$ )
Bob ( $x_2$ )	Lead ( $v_{2,1}$ )
	Dev ( $v_{2,2}$ )
	QA ( $v_{2,3}$ )

EMP	ROLE	EMP	ROLE	EMP	ROLE
Ada	Lead	Ada	Dev	Ada	QA
Bob	Lead	Bob	Lead	Bob	Lead
EMP	ROLE	EMP	ROLE	EMP	ROLE
Ada	Lead	Ada	Dev	Ada	QA
Bob	Dev	Bob	Dev	Bob	Dev
EMP	ROLE	EMP	ROLE	EMP	ROLE
Ada	Lead	Ada	Dev	Ada	QA
Bob	QA	Bob	QA	Bob	QA

*query-answer*



*“Ada and Bob share the same role”*

$$\varphi := (x_1=v_{1,1} \wedge x_2=v_{2,1}) \vee (x_1=v_{1,2} \wedge x_2=v_{2,2}) \vee (x_1=v_{1,3} \wedge x_2=v_{2,3})$$

## (1) Probabilistic Programming Datalog

Bárány, Vince, Balder Ten Cate, Benny Kimelfeld, Dan Olteanu, and Zografoula Vagena. "Declarative probabilistic programming with datalog." *ACM Transactions on Database Systems (TODS)* 42, no. 4 (2017): 1-35.

# Probabilistic Programming Datalog

$\text{weather}(\underline{C}, T, w \in \{\text{sun}, \text{rain}\} \sim \text{Cat}[[P]]) \leftarrow \text{city}(C, P), \text{ts}(T).$

# Probabilistic Programming Datalog

$\text{weather}(\underline{C}, T, w \in \{\text{sun}, \text{rain}\} \sim \text{Cat}[[P]]) \leftarrow \text{city}(C, P), \text{ts}(T).$

`city('Fargo', [.1, .9]), ts('noon')`



# Probabilistic Programming Datalog

$\text{weather}(\underline{C}, T, w \in \{\text{sun}, \text{rain}\} \sim \text{Cat}[[P]]) \leftarrow \text{city}(C, P), \text{ts}(T).$

`city('Fargo', [.1, .9]), ts('noon')`

`weather('Fargo', 'noon', sun) with prob 0.1`

`weather('Fargo', 'noon', rain) with prob 0.9`

# Probabilistic Programming Datalog

$\text{weather}(\underline{C}, T, w \in \{\text{sun}, \text{rain}\} \sim \text{Cat}[[P]]) \leftarrow \text{city}(C, P), \text{ts}(T).$

$\text{city}(\text{'Fargo'}, [.1, .9]), \text{ts}(\text{'noon'})$

$\text{weather}(\text{'Fargo'}, \text{'noon'}, \text{sun})$  with prob 0.1

$\text{weather}(\text{'Fargo'}, \text{'noon'}, \text{rain})$  with prob 0.9

$\text{obs}(\underline{\text{VarId}}, \text{ObsId}, v \in D \sim \text{Cat}[[P]]) \leftarrow \text{lp}(\text{VarId}, D, P),$   
 $\text{sample}(\text{VarId}, \text{ObsId}).$

$\text{lp}(\underline{\text{VarId}}, D, p \in \mathcal{S}_{|D|} \sim \text{Dir}[[H]]) \leftarrow \text{dt}(\text{VarId}, D, H).$

# Where do the constraints come from?

*(Datalog probabilistic program)*

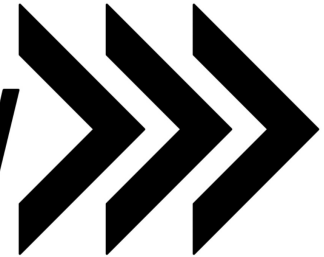


```
dt([red, D], ts, [1, 1, ..., 1]) ← d(D, P, W).  
dt([blue, T], ws, [1, 1, ..., 1]) ← t(T).  
  sample([red, D], P) ← d(D, P, W).  
sample([blue, T], [D, P]) ← d(D, P, W), obs([red, D], P, T).  
  qa*(D, P, W) ← d(D, P, W), obs([blue, T], [D, P], W).
```

APACHE

**ARROW**

***Acero***



*(ground constraints)*



```
(x0,0[s]=v0^x0,1[n]=v0)v(x0,0[s]=v1^x0,1[n]=v1)  
(x0,0[e]=v0^x1,0[w]=v0)v(x0,0[e]=v1^x0,1[w]=v1)  
(x1,0[s]=v0^x1,1[n]=v0)v(x1,0[s]=v1^x1,1[n]=v1)  
(x0,1[e]=v0^x1,1[w]=v0)v(x0,1[e]=v1^x1,1[w]=v1)  
(x2,0[s]=v0^x2,1[n]=v0)v(x2,0[s]=v1^x2,1[n]=v1)
```

...

...

(many)



**LLVM/ClangJIT**

# Inference in Action



(a) Noisy Image



(b) Gibbs sampler  
denoising



(c) Variational inference  
denoising

**Thank you!**

**(Questions?)**