

# Thinking Compositionally about Inference

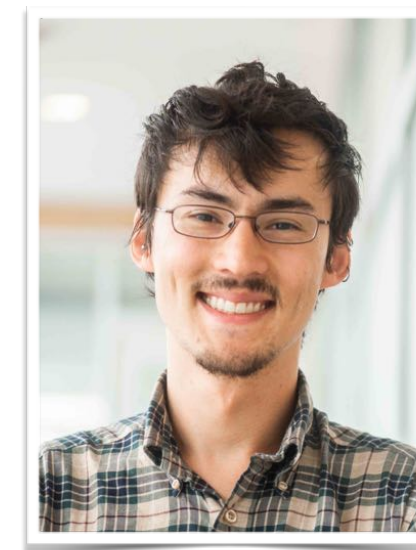
Jan-Willem van de Meent



Heiko Zimmermann



Babak Esmaeili



Sam Stites

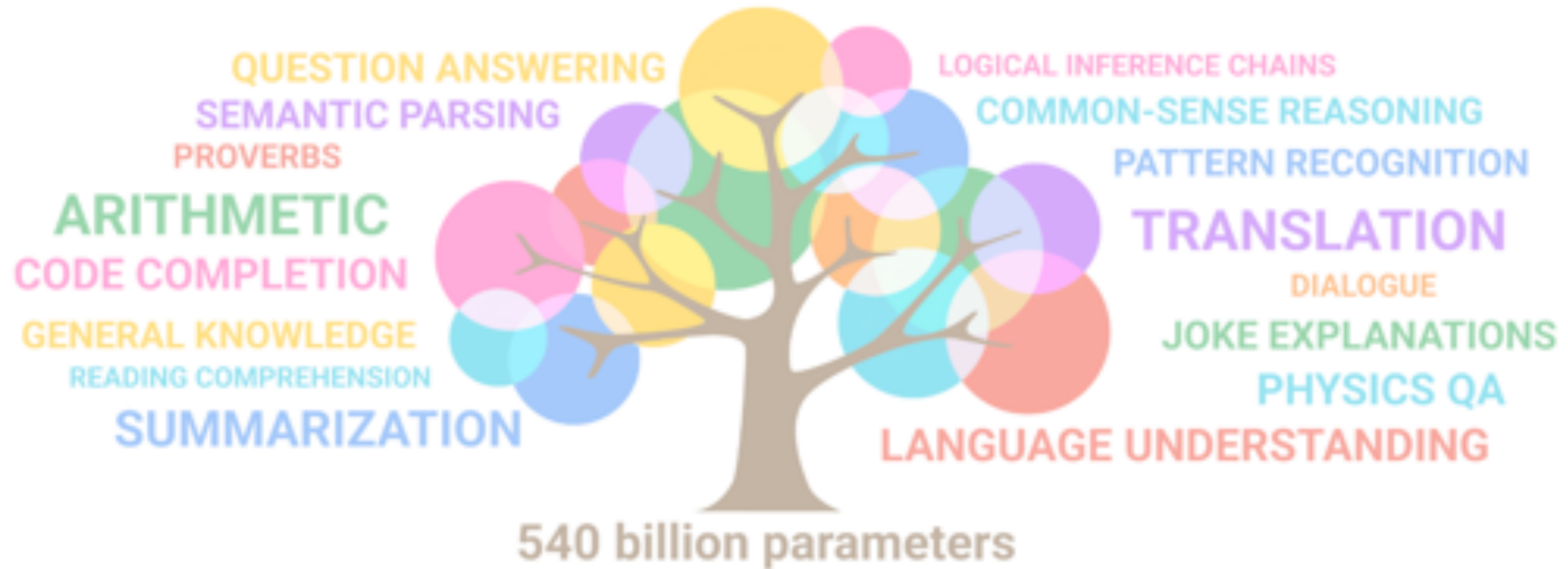


Hao Wu



Eli Sennesh

# Deep Learning Success Stories



PaLM

# Deep Learning Success Stories

## Explaining a joke

### Prompt

Explain this joke:

Joke: Did you see that Google just hired an eloquent whale for their TPU team? It showed them how to communicate between two different pods!

### Model Response

Prediction: TPUs are a type of computer chip that Google uses for deep learning. A "pod" is a group of TPUs. A "pod" is also a group of whales. The joke is that the whale is able to communicate between two groups of whales, but the speaker is pretending that the whale is able to communicate between two groups of TPUs.

PaLM

# Deep Learning Success Stories



Painting of the orange cat Otto von Garfield, Count of Bismarck-Schönhausen, Duke of Lauenburg, Minister-President of Prussia. Depicted wearing a Prussian Pickelhaube and eating his favorite meal - lasagna.



A photo of the back of a wombat wearing a **backpack** and holding a **walking stick**. It is next to a **waterfall** and is staring at a **distant mountain**.

Dall-E 2, Imagen, Parti

# Is Scale All We Need?

## **The Bitter Lesson**

**Rich Sutton, March 13, 2019**

The biggest lesson that can be read from 70 years of AI research is that general methods that leverage computation are ultimately the most effective, and by a large margin.

<http://www.incompleteideas.net/Incldeas/BitterLesson.html>

- Scale can lead to abstractions and generalization across tasks
- Still difficult to know when a model will succeed or fail.
- How can we scale up to more diverse application domains?

# Is Scale All We Need?

A superintelligent chess AI with 5000 ELO is playing a game of chess against a human. The AI is playing as black. This is a transcript of the game.

1. e4 e5
2. Nf3 Nc6
3. Bb5 a6
4. Bxc6 dxc6
5. O-O Qf6
6. d3 Qg6
7. Nxe5 Qxe4
8. dxe4 Bd6
9. Bf4 Bxe5
10. Bxe5 Ne7
11. Bxc7 Nxc6



- Scale can lead to abstractions and generalization across tasks
- Still difficult to know when a model will succeed or fail.
- How can we scale up to more diverse application domains?

<https://jacobbuckman.com/2022-06-14-an-actually-good-argument-against-naive-ai-scaling/>

# Adapting Deep Learning to New Domains

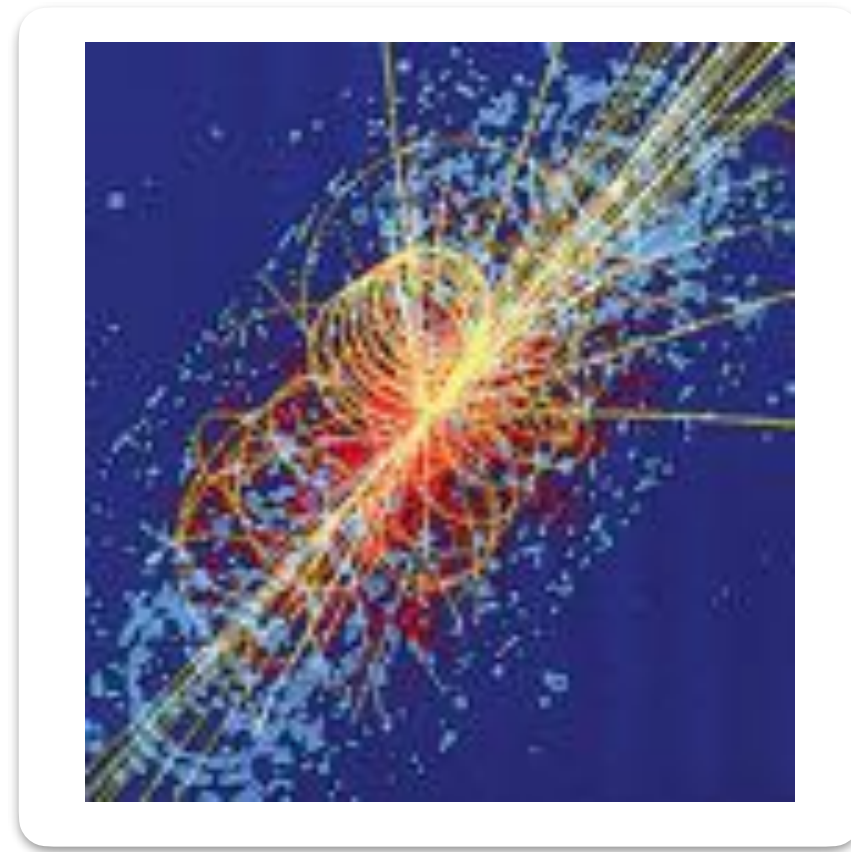


## Lessons from deep learning

1. Gradient descent scales *really* well
2. Model engineering scales *pretty* well

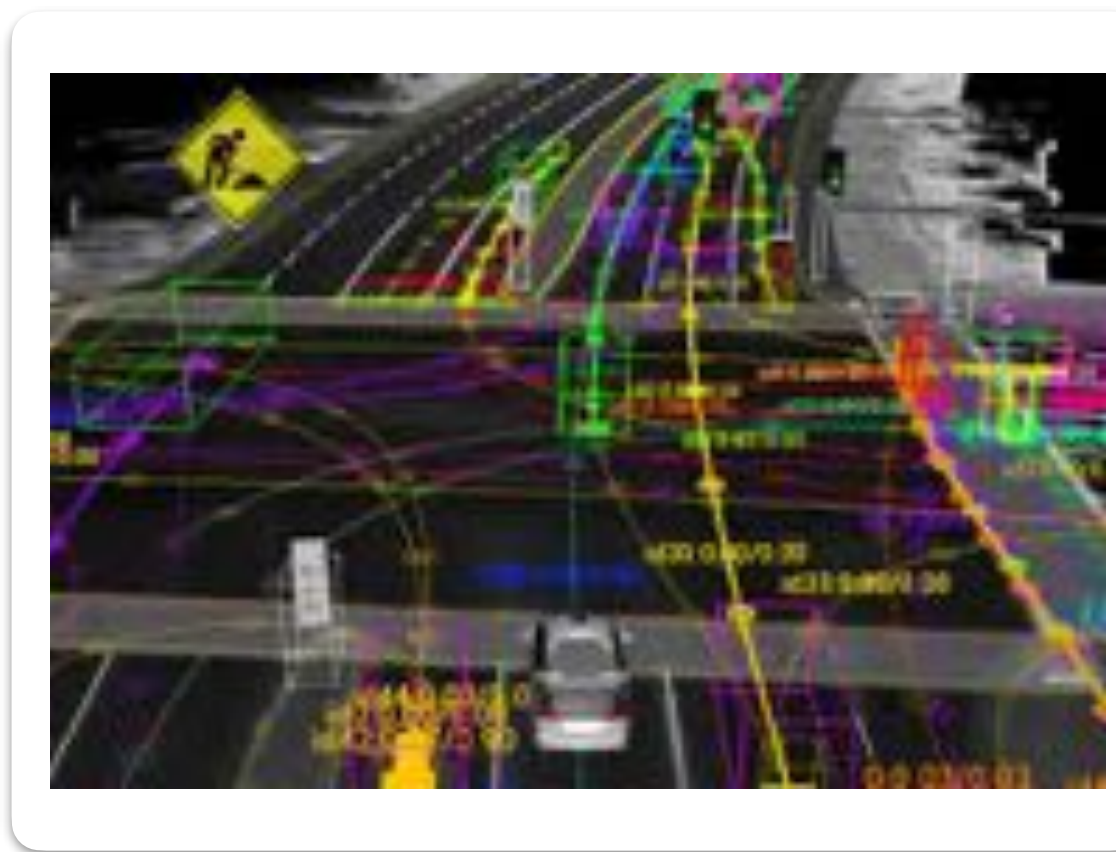
# Horizons of AI Research

## Science & Engineering



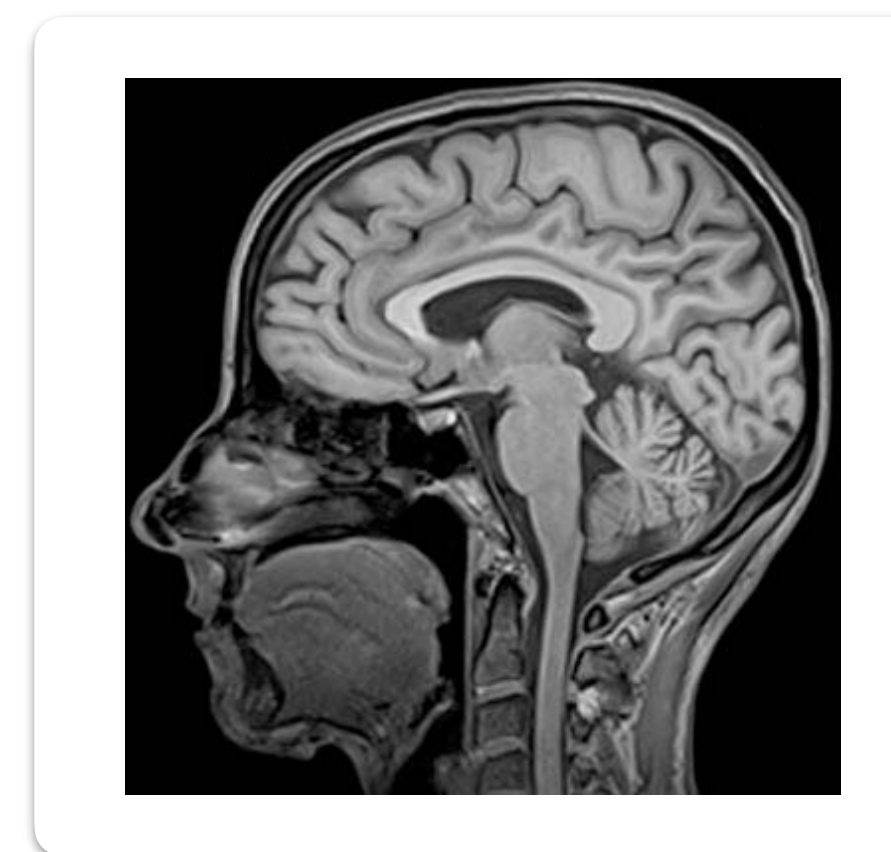
Deep domain knowledge  
but limited data

## Autonomous Vehicles



Generalization to  
long tail events

## Healthcare



Many prediction tasks,  
imbalanced data

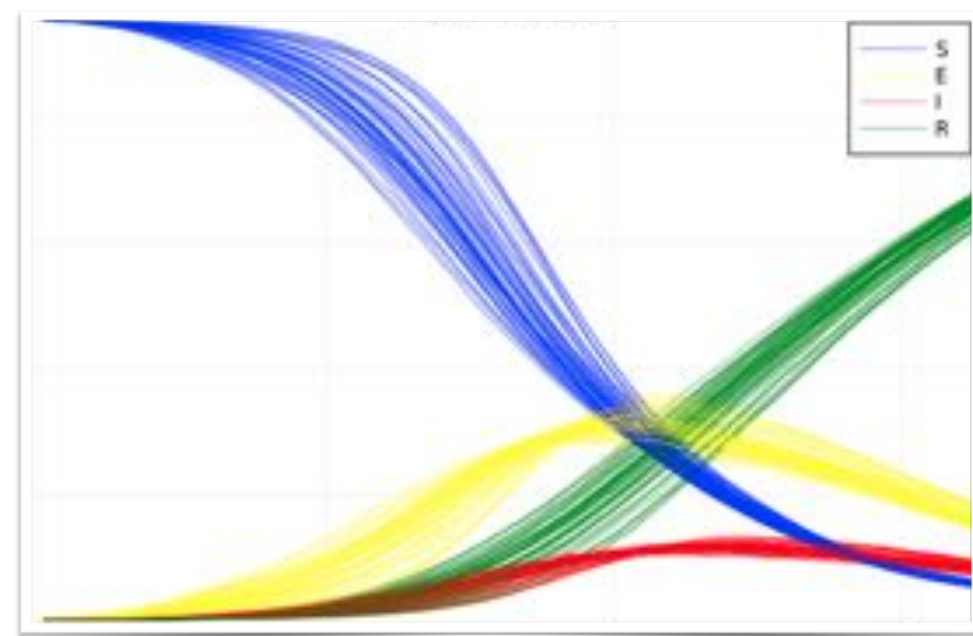
## Challenges in emerging domains

1. Incorporating (enough) domain knowledge
2. Reliable generalization across related tasks
3. Avoiding overconfident predictions



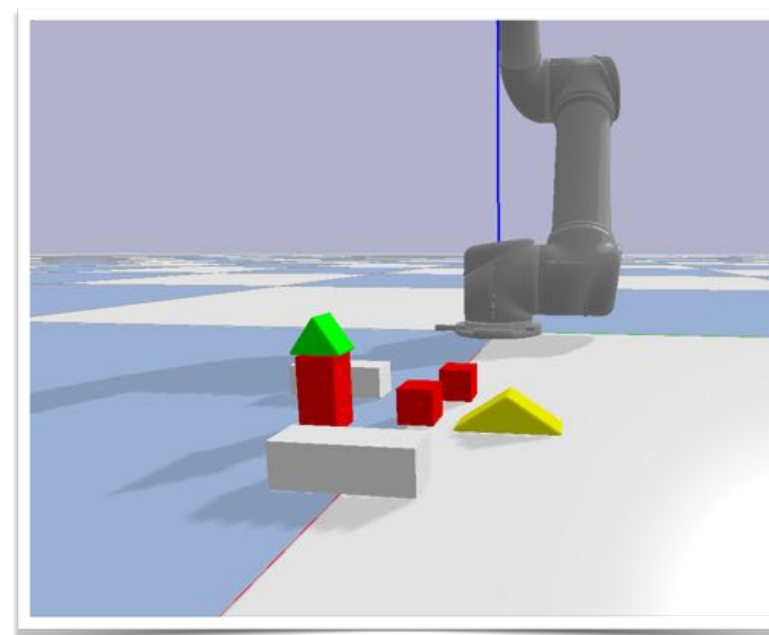
# What Models are Useful?

## Simulation-based Modeling



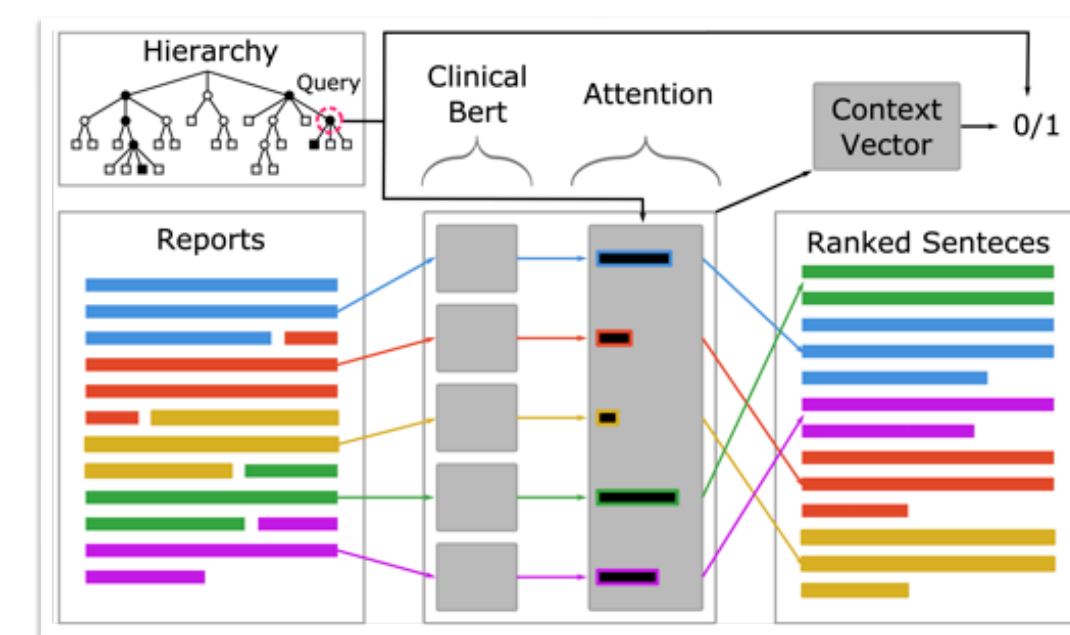
[Smedemark-Margulies et al., 2021]

## Planning and Robotics



[Biza et al., 2021]

## Vision & Language



[McInerney et al., 2020]

*Stronger assumptions*

*Weaker assumptions*



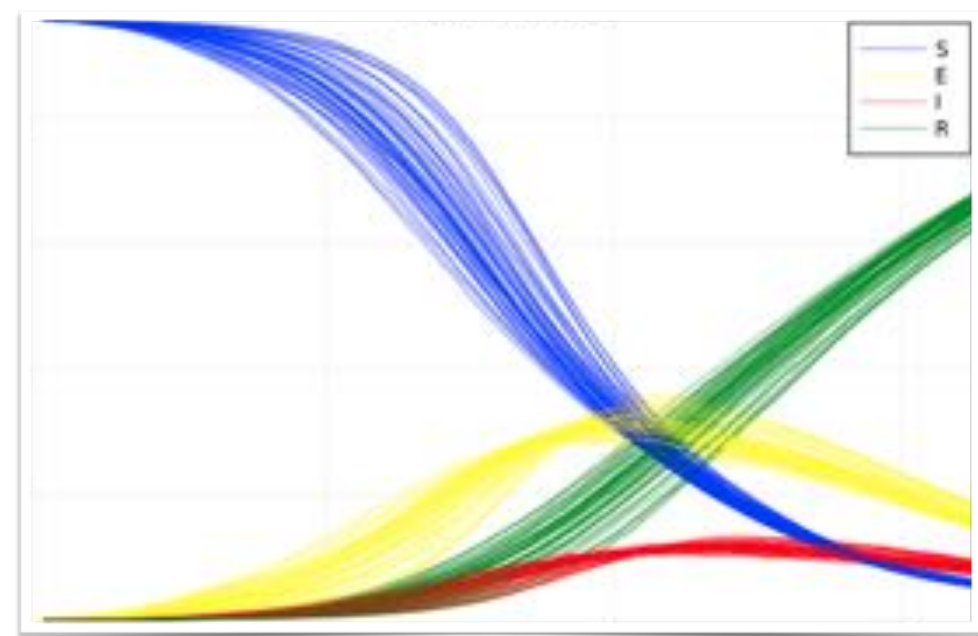
Known dynamics  
(e.g. PDEs) for system

More knowledge  
(and edge cases)

Some domain knowledge  
(e.g. structure)

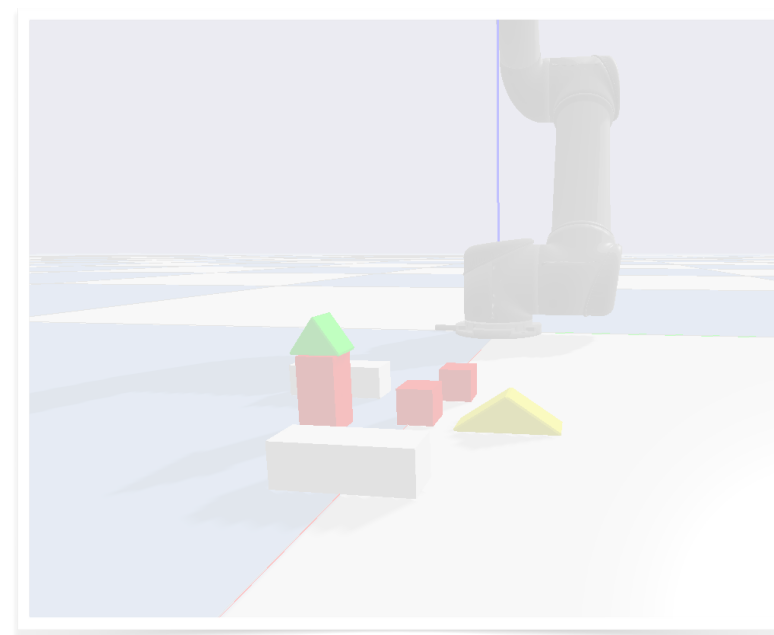
# What Models are Useful?

## Simulation-based Modeling



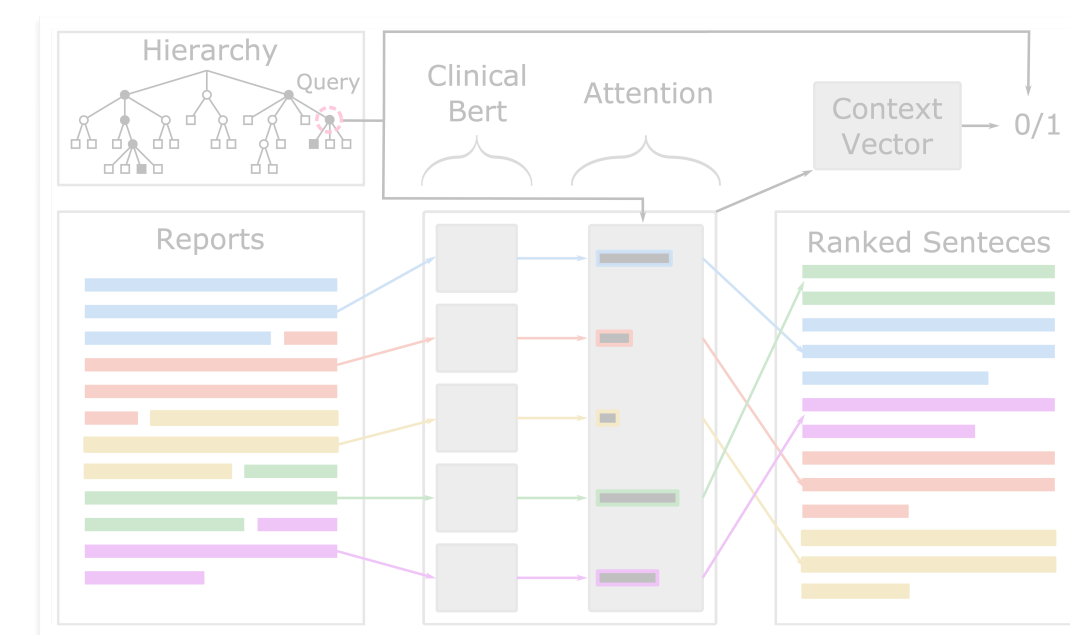
[Smedemark-Margulies et al., 2021]

## Planning and Robotics



[Biza et al., 2021]

## Vision & Language



[McInerney et al., 2020]

*Stronger assumptions*

*Weaker assumptions*



Known dynamics  
(e.g. PDEs) for system

More knowledge  
(and edge cases)

Some domain knowledge  
(e.g. structure)

# The Next 700 AI Domains

# The Next 700 Programming Languages

P. J. Landin

*Univac Division of Sperry Rand Corp., New York, New York*

“... today ... 1,700 special programming languages used to ‘communicate’ in over 700 application areas.”—*Computer Software Issues*, an American Mathematical Association Prospectus, July 1965.

**Volume 9 / Number 3 / March, 1966**

**Communications of the ACM**

**157**

## Two Ingredients for a Language

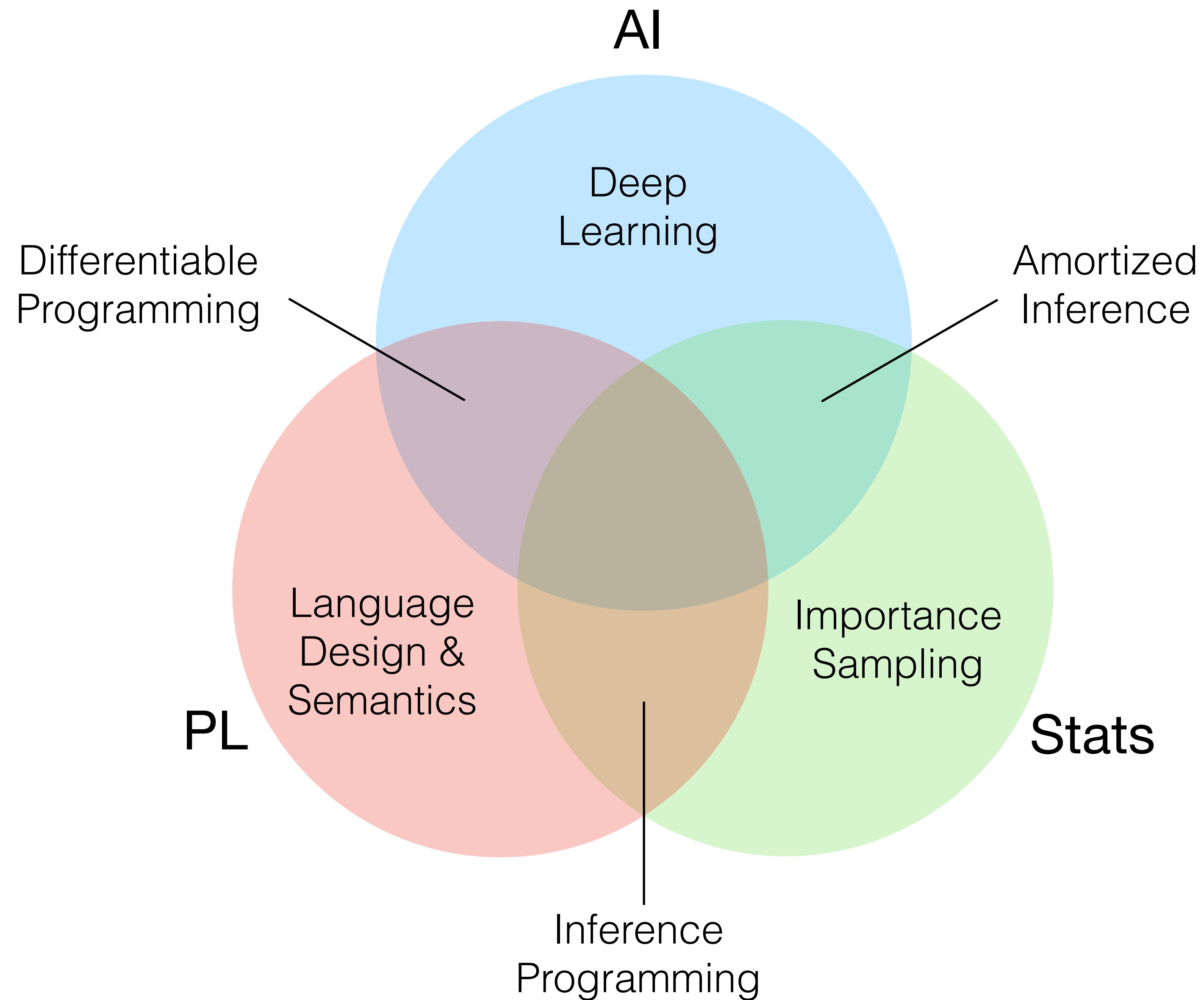
1. Core operations / abstractions
2. Mechanisms for composition into program

# Differentiable Programming

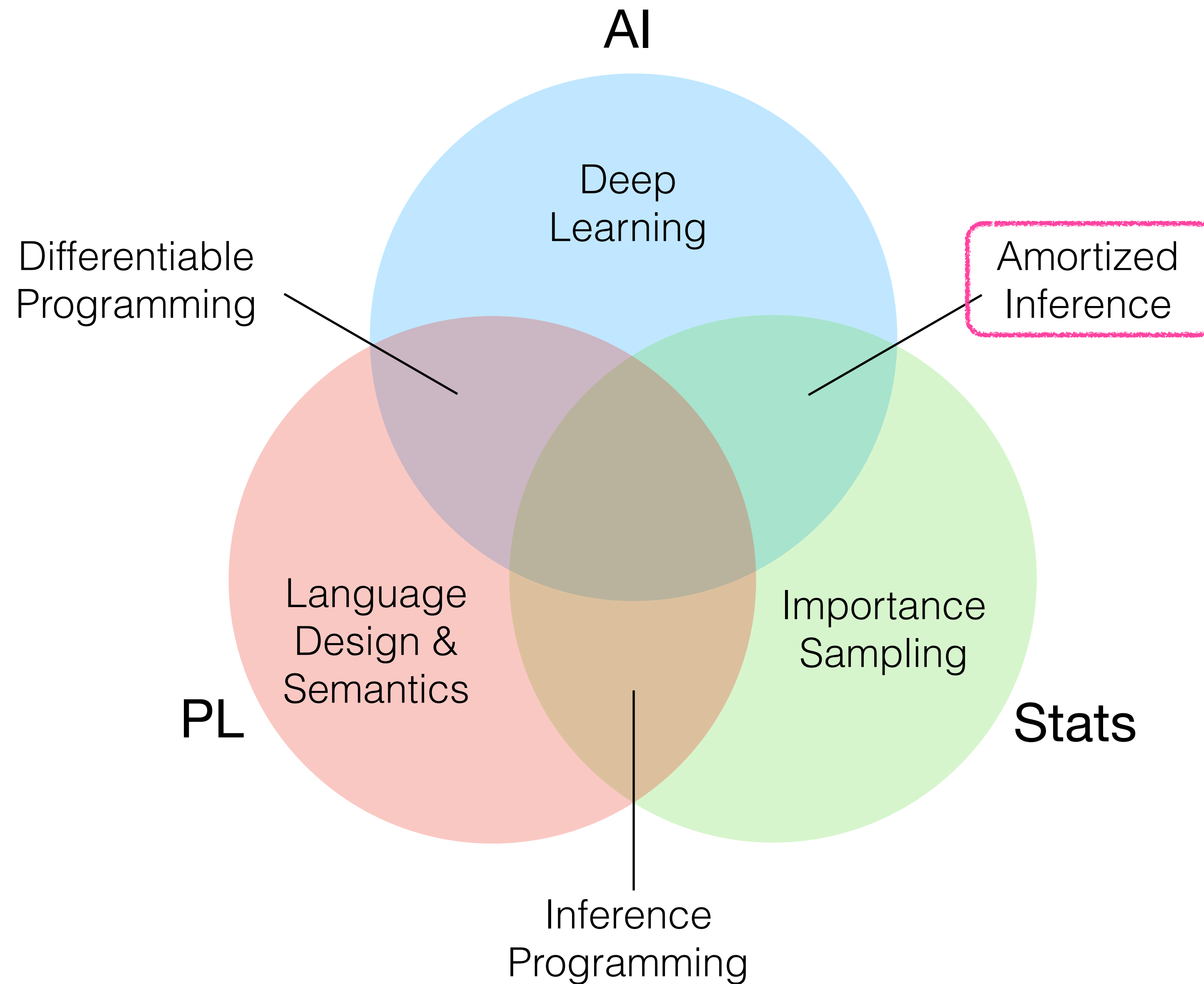


1. Abstractions: Differentiation, Tensor Calculus, Layers
2. Composition: Networks, Objectives, Optimization

# Deep Probabilistic Programming

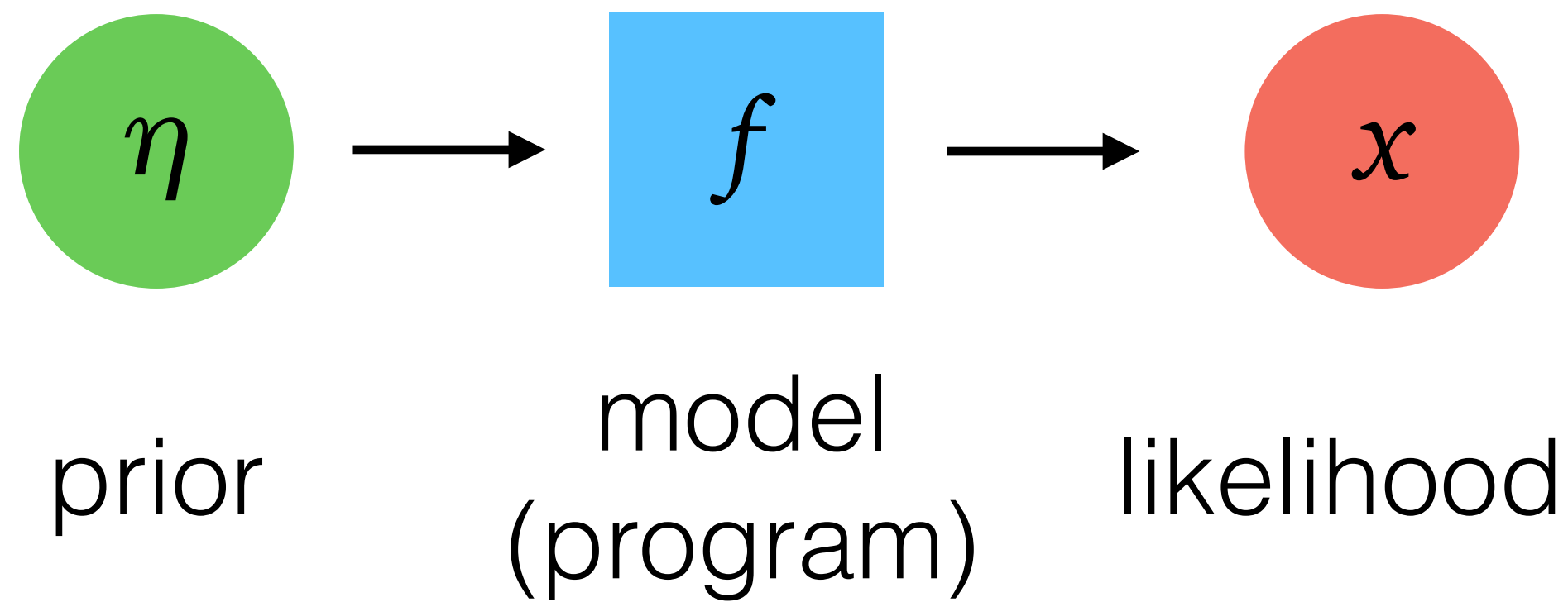


# Deep Probabilistic Programming

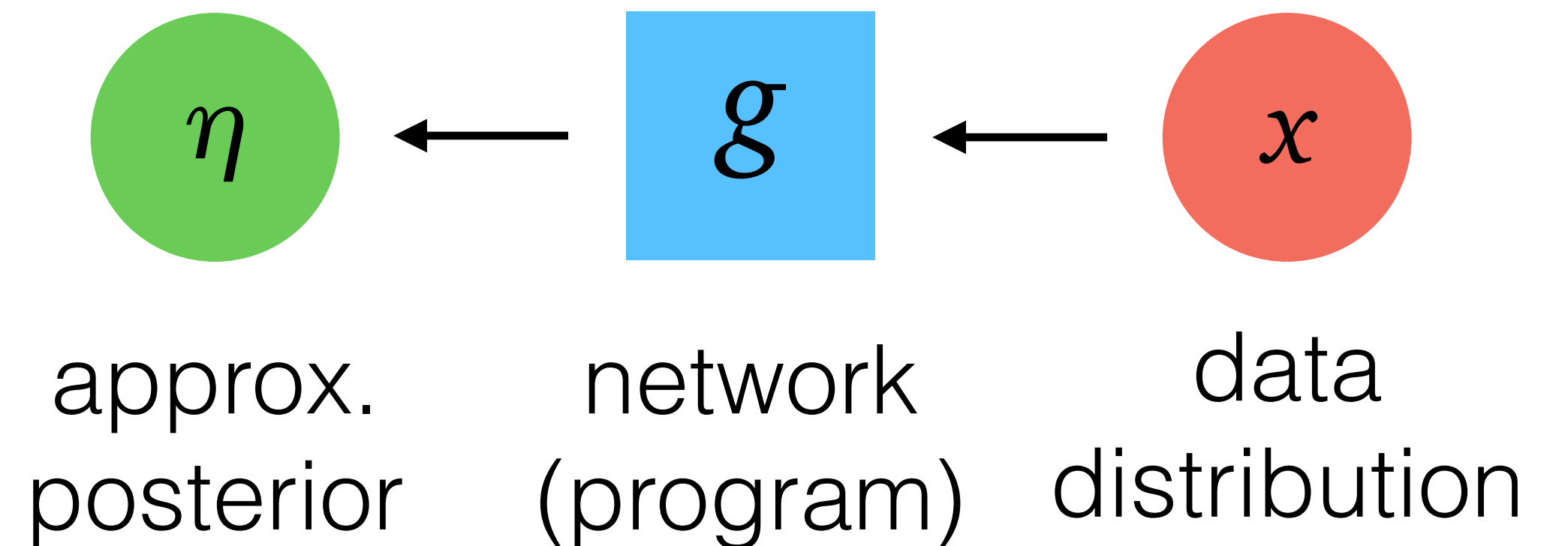


# Amortized Inference

Stochastic Simulator  
*(most of science and engineering)*



Inference Model  
*(approximate inverse)*



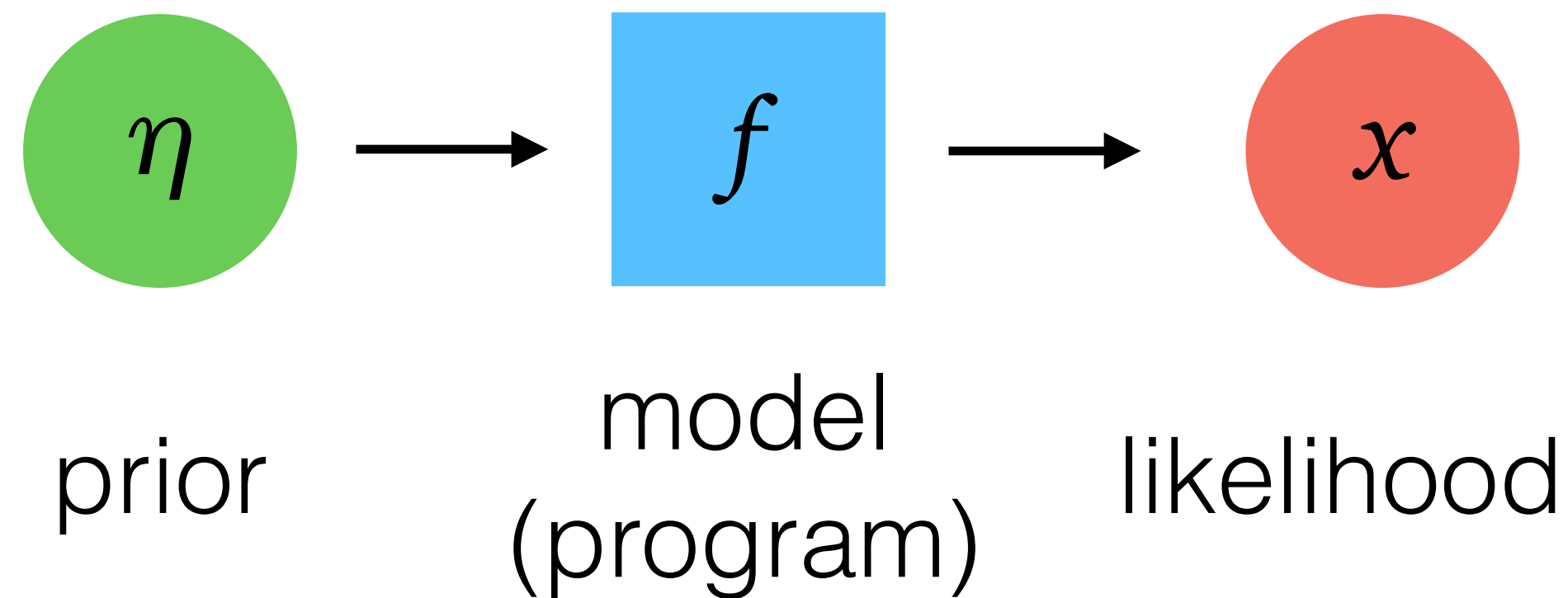
$$p_{\theta}(x, \eta) = p(\eta) p(x | f_{\theta}(\eta))$$

$$p_{\theta}(\eta | x) = \frac{p(\eta) p_{\theta}(x | \eta)}{p_{\theta}(x)}$$

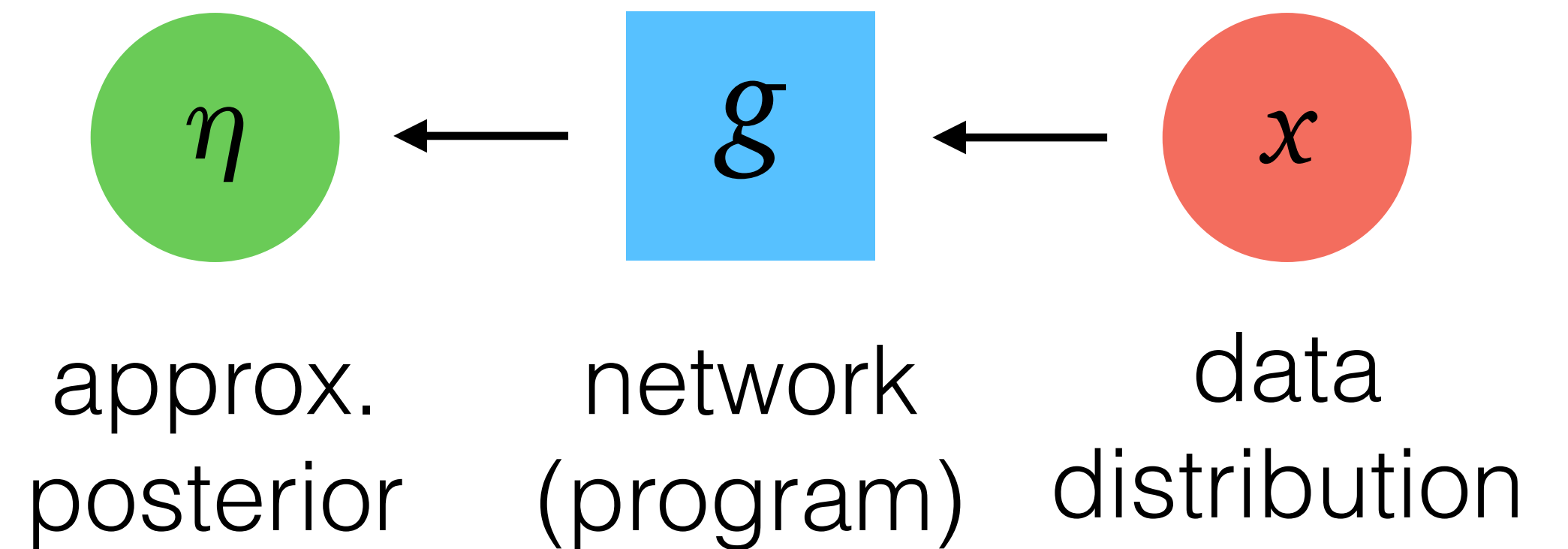


# Amortized Inference

Stochastic Simulator  
*(most of science and engineering)*



Inference Model  
*(approximate inverse)*



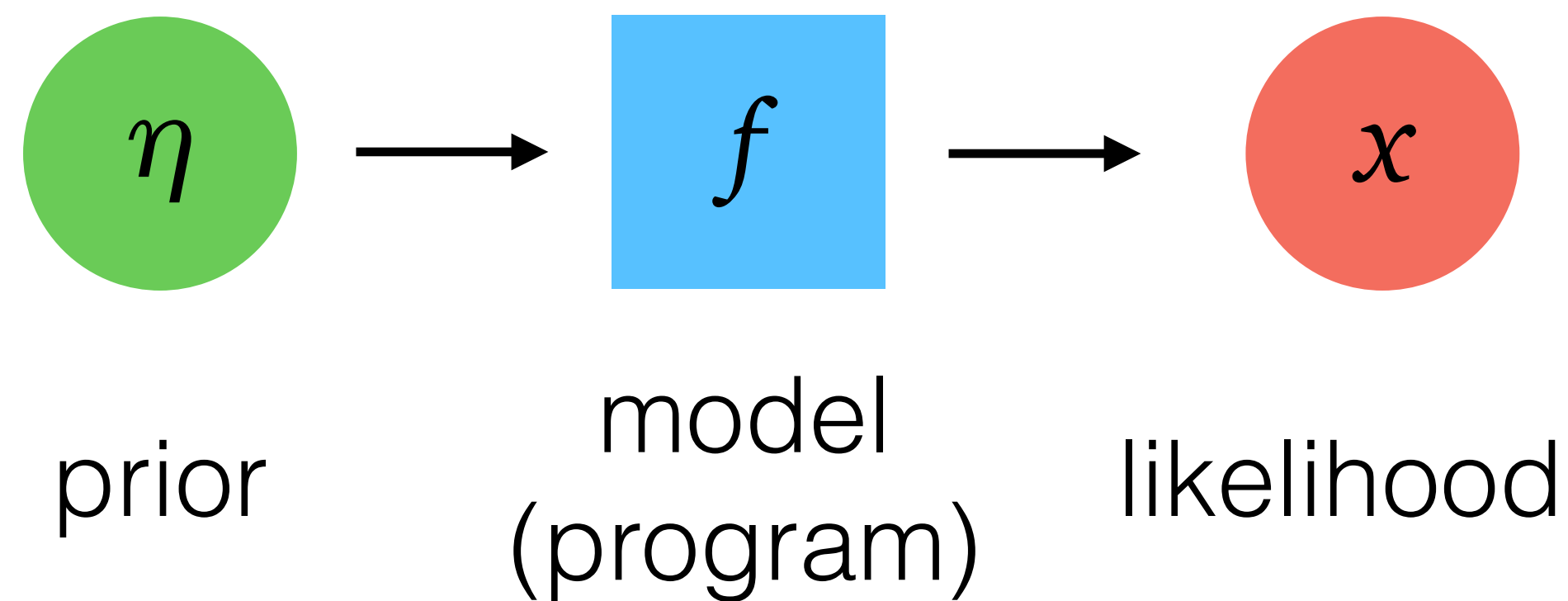
$$p_{\theta}(x, \eta) = p(\eta) p(x | f_{\theta}(\eta))$$

$$p_{\theta}(\eta | x) = \frac{p_{\theta}(x | \eta) p(\eta)}{p_{\theta}(x)}$$

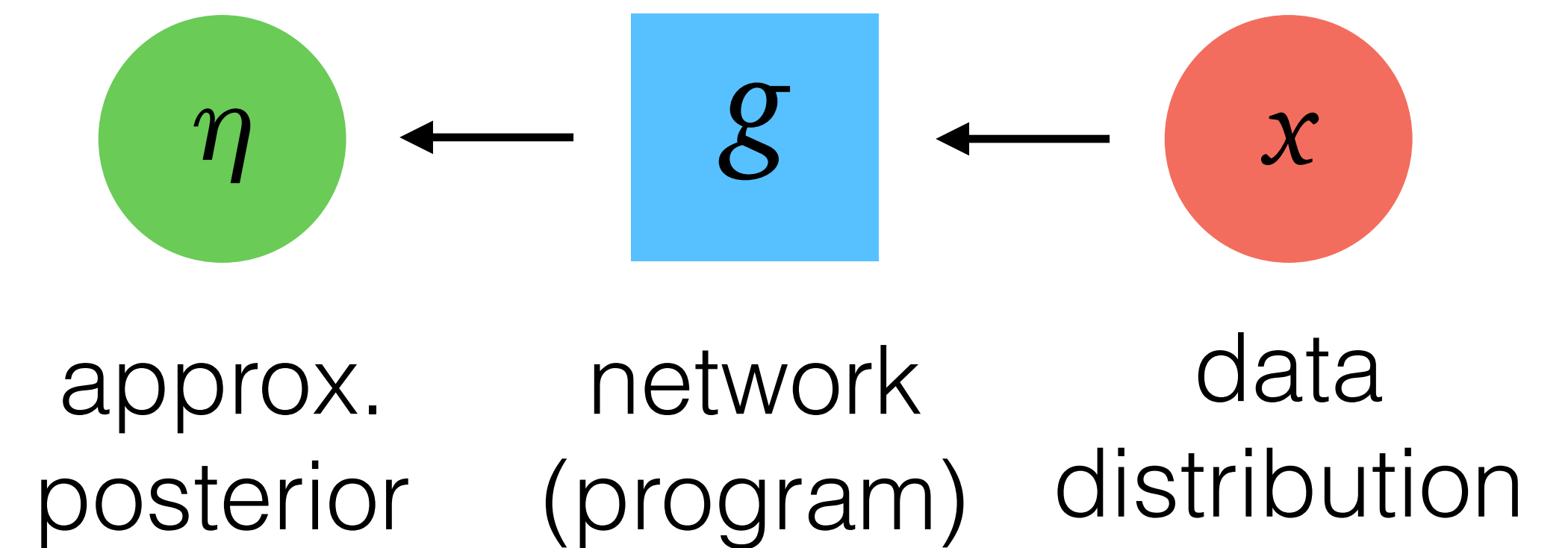
$$q_{\phi}(x, \eta) = q(x) q(\eta | g_{\phi}(x))$$

# Amortized Inference

Generative Model  
(*stochastic simulator*)



Inference Model  
(*approximate inverse*)



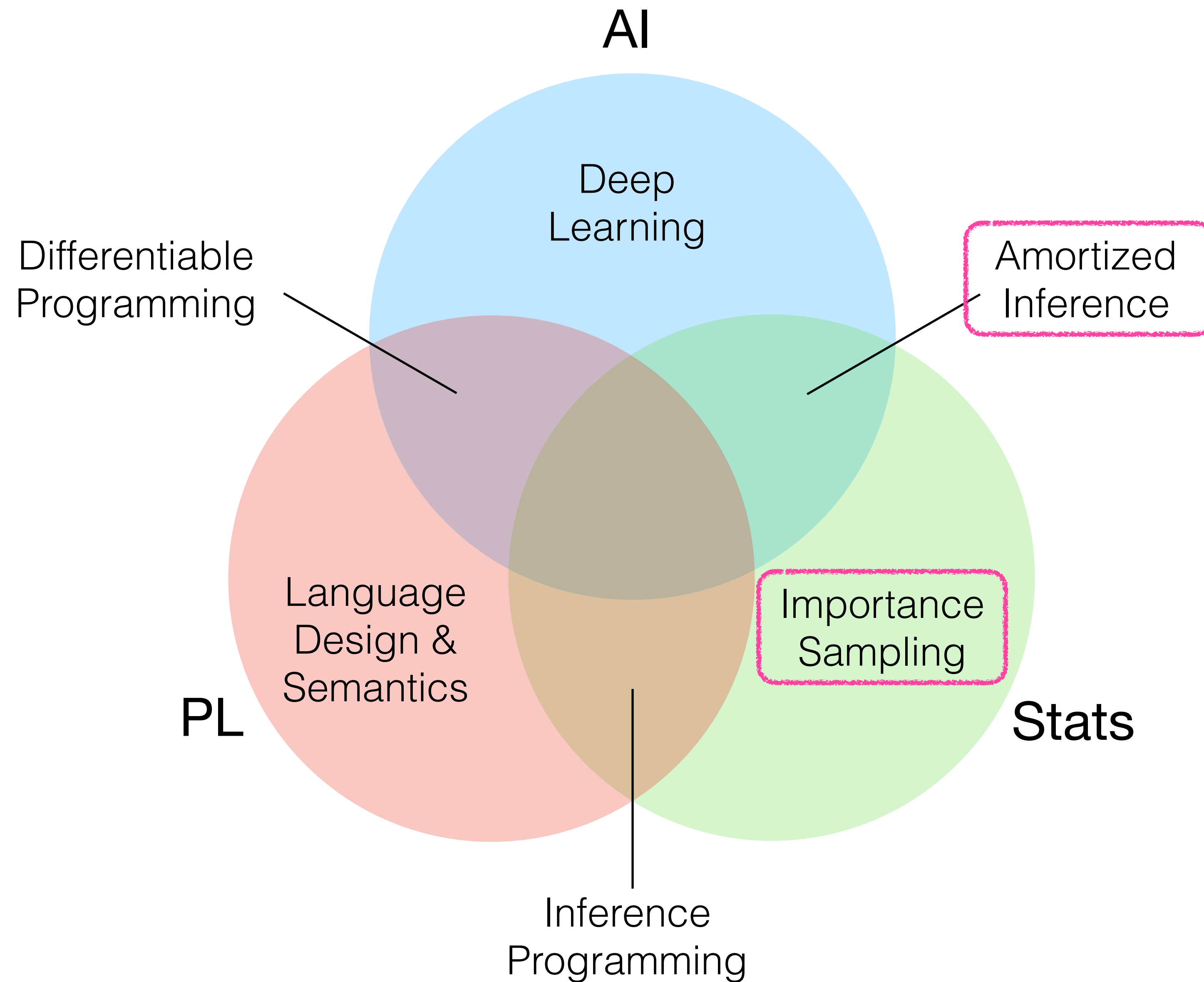
Model Learning

$$\min_{\theta} D(q(x) \parallel p_{\theta}(x))$$

Amortized Inference

$$\min_{\phi} \mathbb{E}_{x \sim q} [D(q_{\phi}(\eta | x) \parallel p_{\theta}(\eta | x))]$$

# Deep Probabilistic Programming



# Minimizing the Inclusive KL divergence

Idea 1: Minimize *inclusive* KL (rather than *exclusive* KL)

$$\min_{\phi} \mathbb{E}_{x \sim q} [D_{\text{KL}}(q_{\phi}(\eta | x) || p_{\theta}(\eta | x))] \rightarrow \min_{\phi} \mathbb{E}_{x \sim q} [D_{\text{KL}}(p_{\theta}(\eta | x) || q_{\phi}(\eta | x))]$$

Idea 2: Use importance sampling to approximate gradient

$$-\nabla_{\phi} D_{\text{KL}}(p_{\theta}(\eta | x) || q_{\phi}(\eta | x)) = \mathbb{E}_{\eta \sim p_{\theta}(\cdot | x)} [\nabla_{\phi} \log q_{\phi}(\eta | x)]$$

Use importance sampling

$$w^l = \frac{p_{\theta}(x, \eta^l)}{q_{\phi}(\eta^l | x)} \quad \eta^l \sim q_{\phi}(\eta | x)$$

# Minimizing the Inclusive KL divergence

Idea 1: Minimize *inclusive* KL rather than *exclusive* KL

$$\min_{\phi} \mathbb{E}_{x \sim q} \left[ D_{\text{KL}}(q_{\phi}(\eta | x) || p_{\theta}(\eta | x)) \right] \rightarrow \min_{\phi} \mathbb{E}_{x \sim q} \left[ D_{\text{KL}}(p_{\theta}(\eta | x) || q_{\phi}(\eta | x)) \right]$$

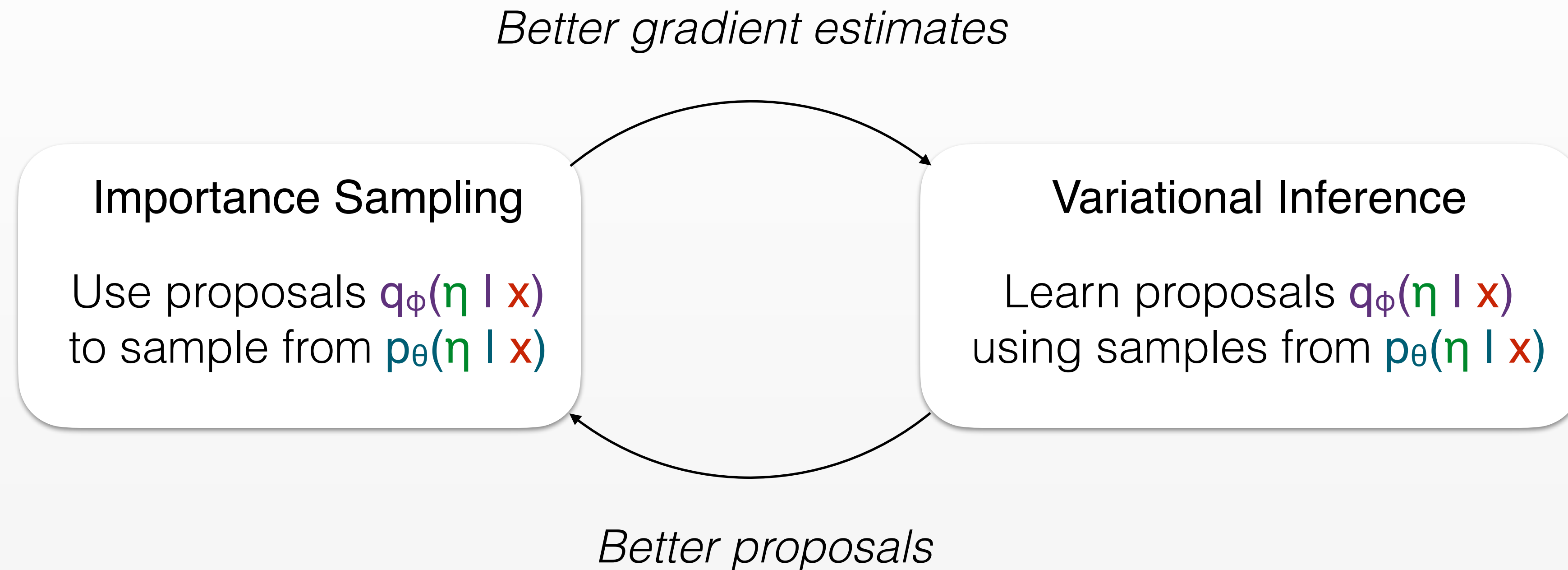
Idea 2: Use importance sampling to approximate gradient

$$-\nabla_{\phi} D_{\text{KL}}(p_{\theta}(\eta | x) || q_{\phi}(\eta | x)) = \mathbb{E}_{\eta \sim p_{\theta}(\cdot | x)} \left[ \nabla_{\phi} \log q_{\phi}(\eta | x) \right]$$

Use importance sampling

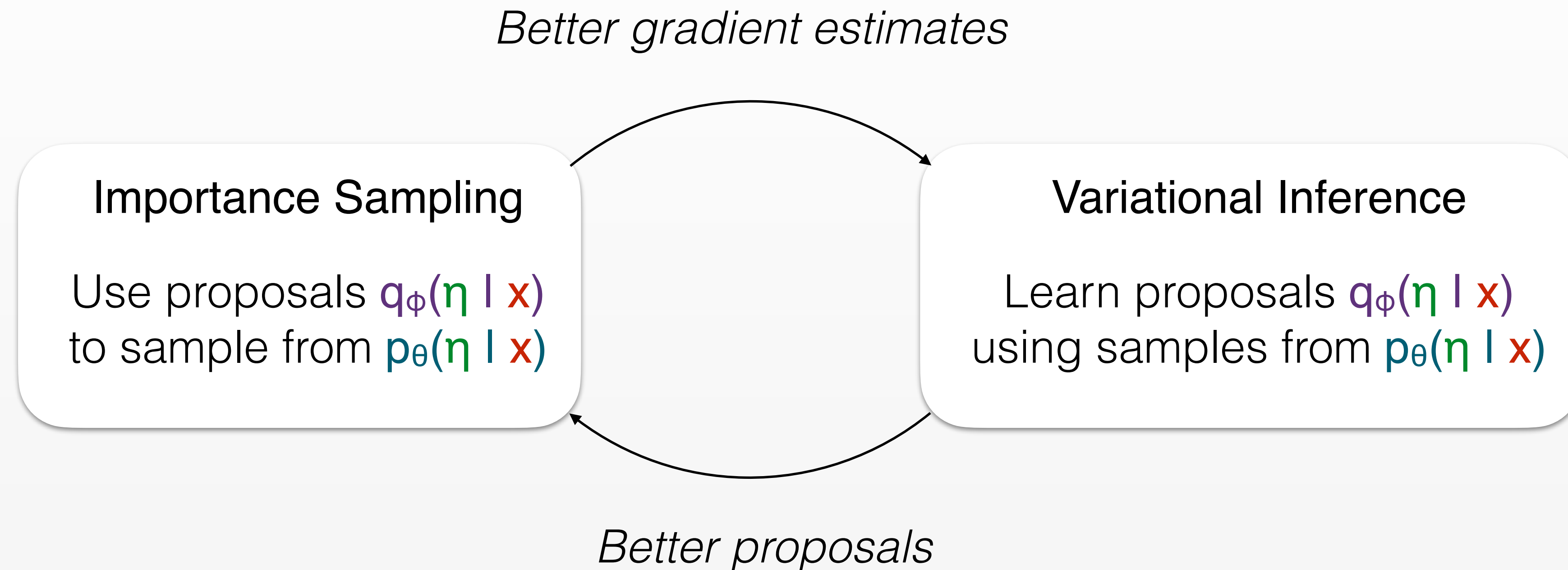
$$w^l = \frac{p_{\theta}(x, \eta^l)}{q_{\phi}(\eta^l | x)} \quad \eta^l \sim q_{\phi}(\eta | x) \quad \approx \sum_{l=1}^L \frac{w^l}{\sum_{l'} w^{l'}} \nabla_{\phi} \log q_{\phi}(\eta^l | x)$$

# Amortized Importance Samplers



- Does not rely on differentiable models / reparameterization
- Often works as well as, or better than, maximizing a lower bound

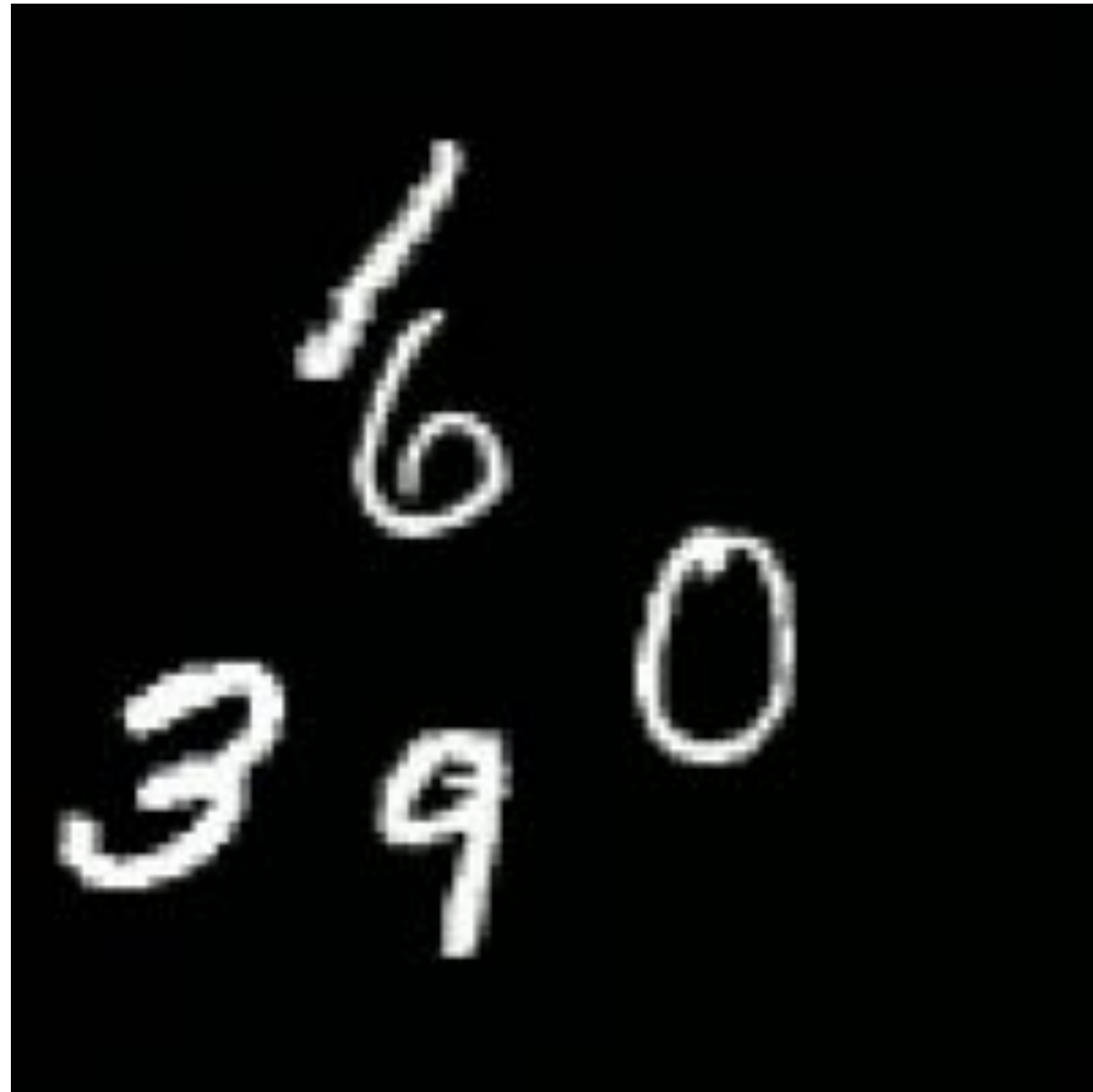
# Amortized Importance Samplers



**Opportunity:** New VI methods based on  
SMC samplers, nested importance samplers, etc

# Example (~2019): Amortized Population Gibbs

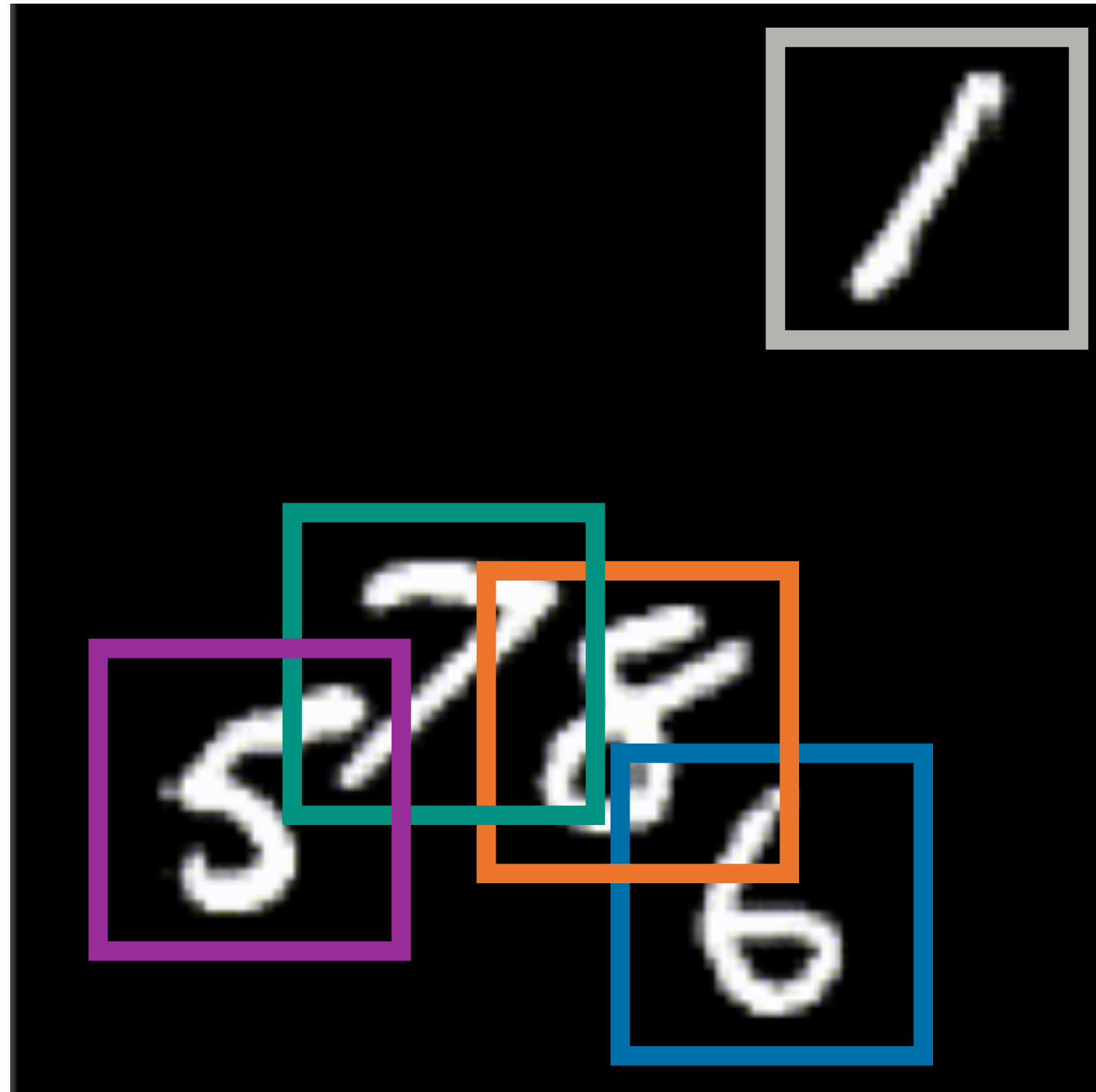
## Task: Unsupervised Tracking



- Corpus level (*many videos*)  
Digit shapes  
Transition dynamics
- Instances (*single videos*)  
Object representations
- Data-points (*single frames*)  
Object positions



# Example (~2019): Amortized Population Gibbs



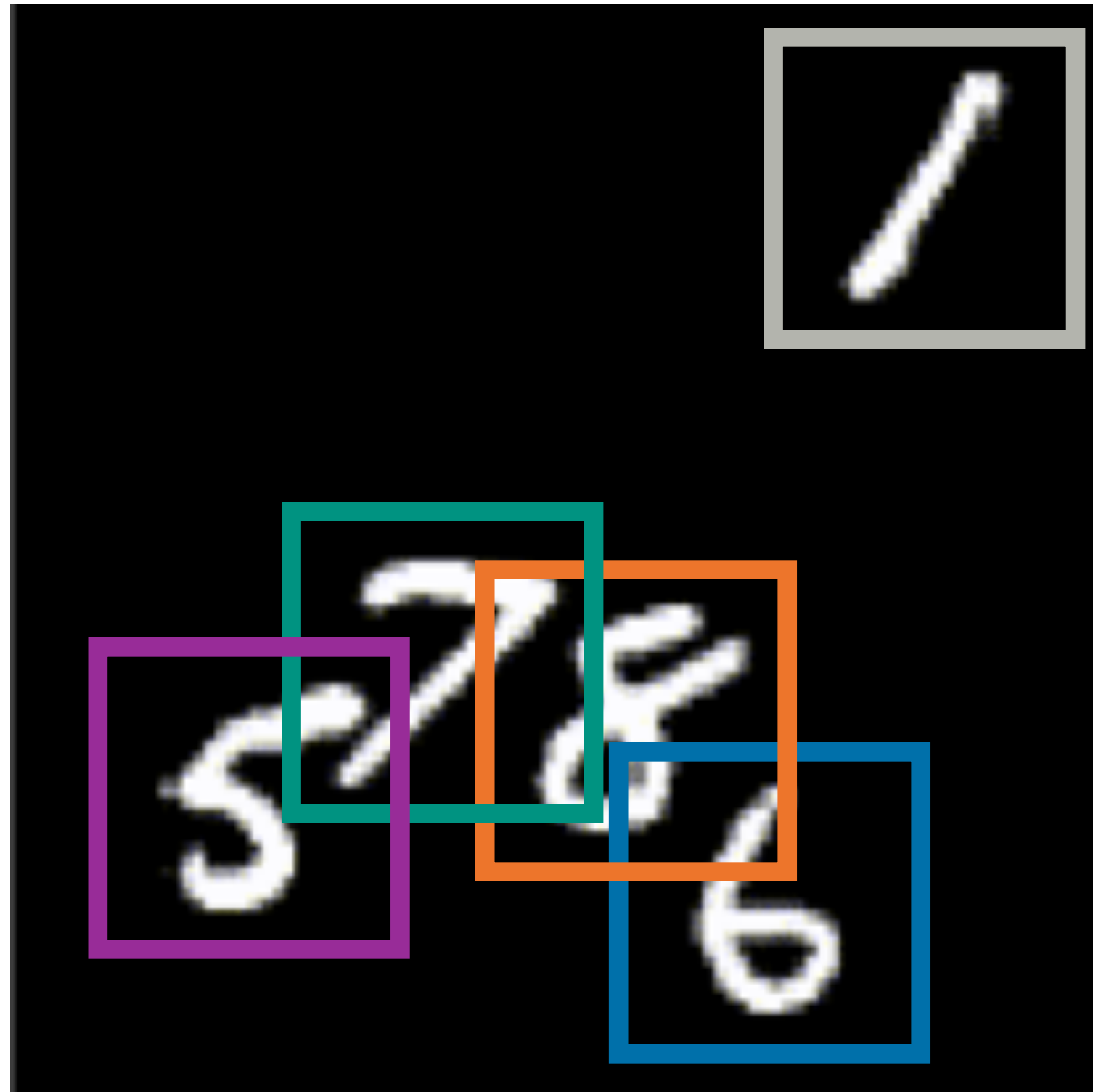
## Classic Chicken-and-Egg Problem

- **Easy:** Infer object representations given object positions



- **Also Easy:** Infer positions given object representations
- **Not Easy:** Joint inference of positions and representations

# Example (~2019): Amortized Population Gibbs



## Classic Solution: Iterate

- **Step 0:** Initialize representations and positions.
- **Update 1:** Infer object representations given object positions
- **Update 2:** Infer object representations given object positions

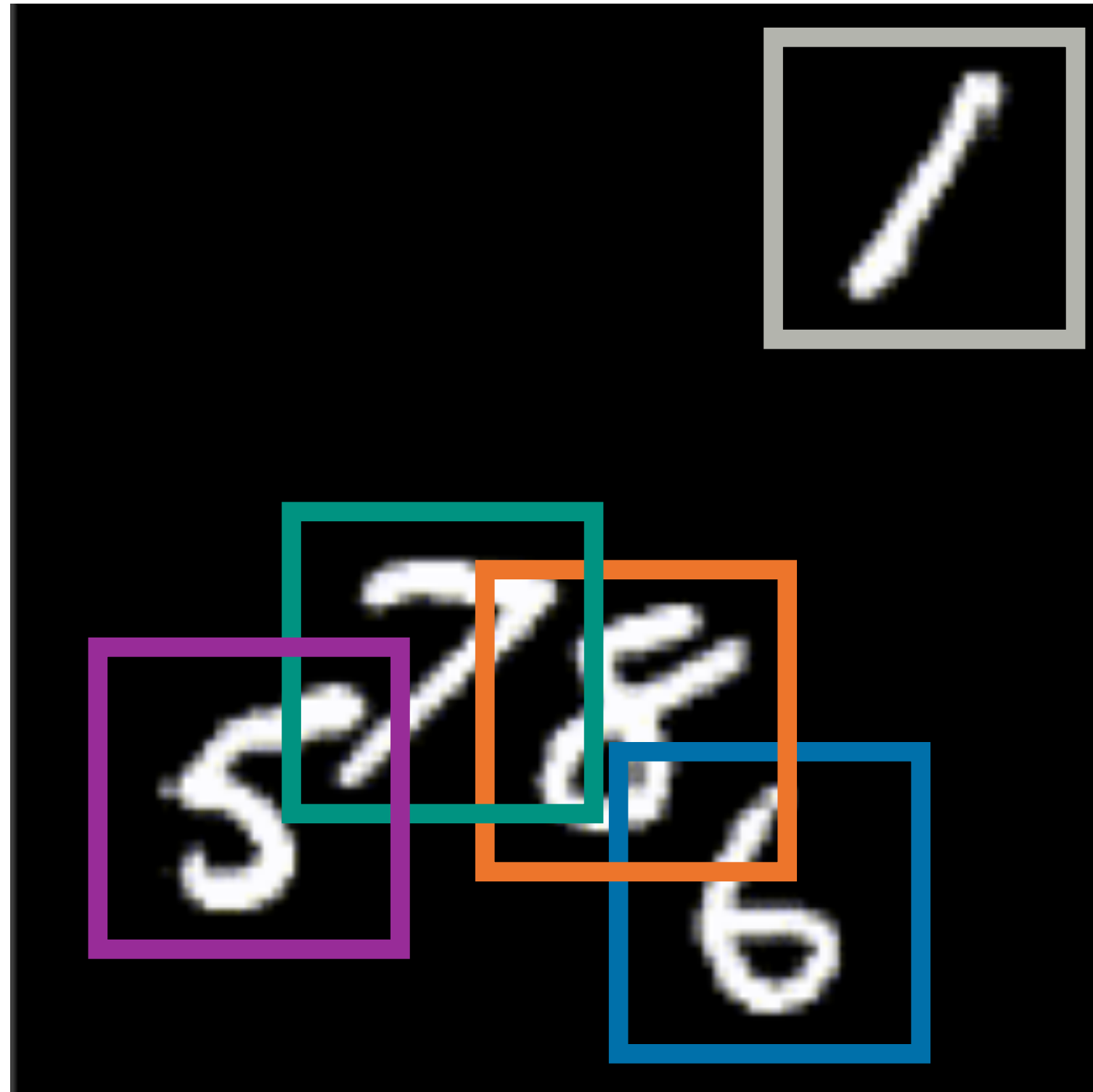
$$\eta \sim p(\eta | x, z)$$

$$z \sim p(z | x, \eta)$$

## Problem:

Only computable in conjugate exponential family models

# Example (~2019): Amortized Population Gibbs



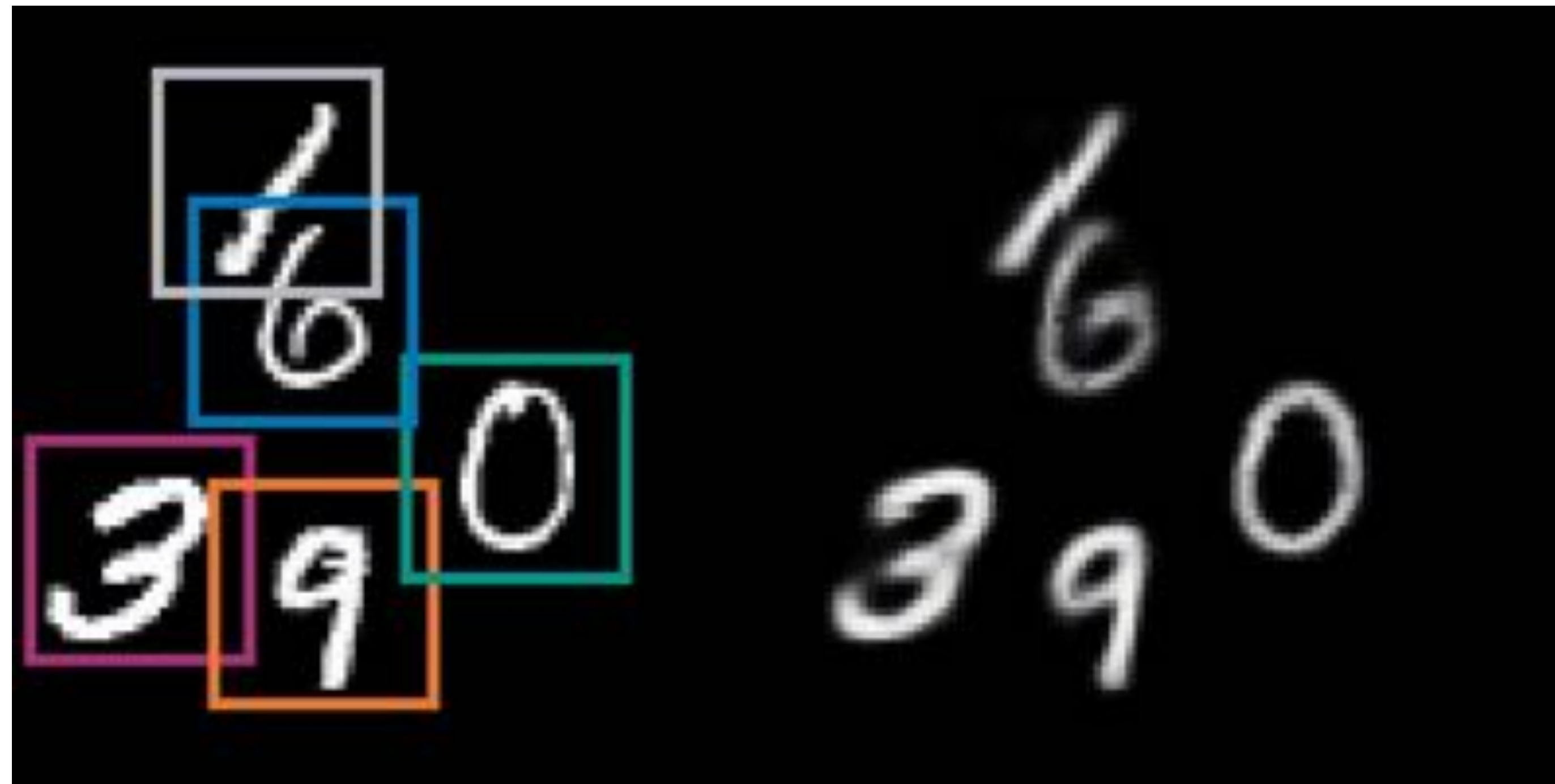
Modern solution: **Learn** Updates

- **Step 0:** Initialize representations and positions.
- **Update 1:** Infer object representations given object positions  
 $\eta \sim q_{\phi}(\eta | x, z)$
- **Update 2:** Infer object representations given object positions  
 $z \sim q_{\phi}(z | x, \eta)$

# Example (~2019): Amortized Population Gibbs

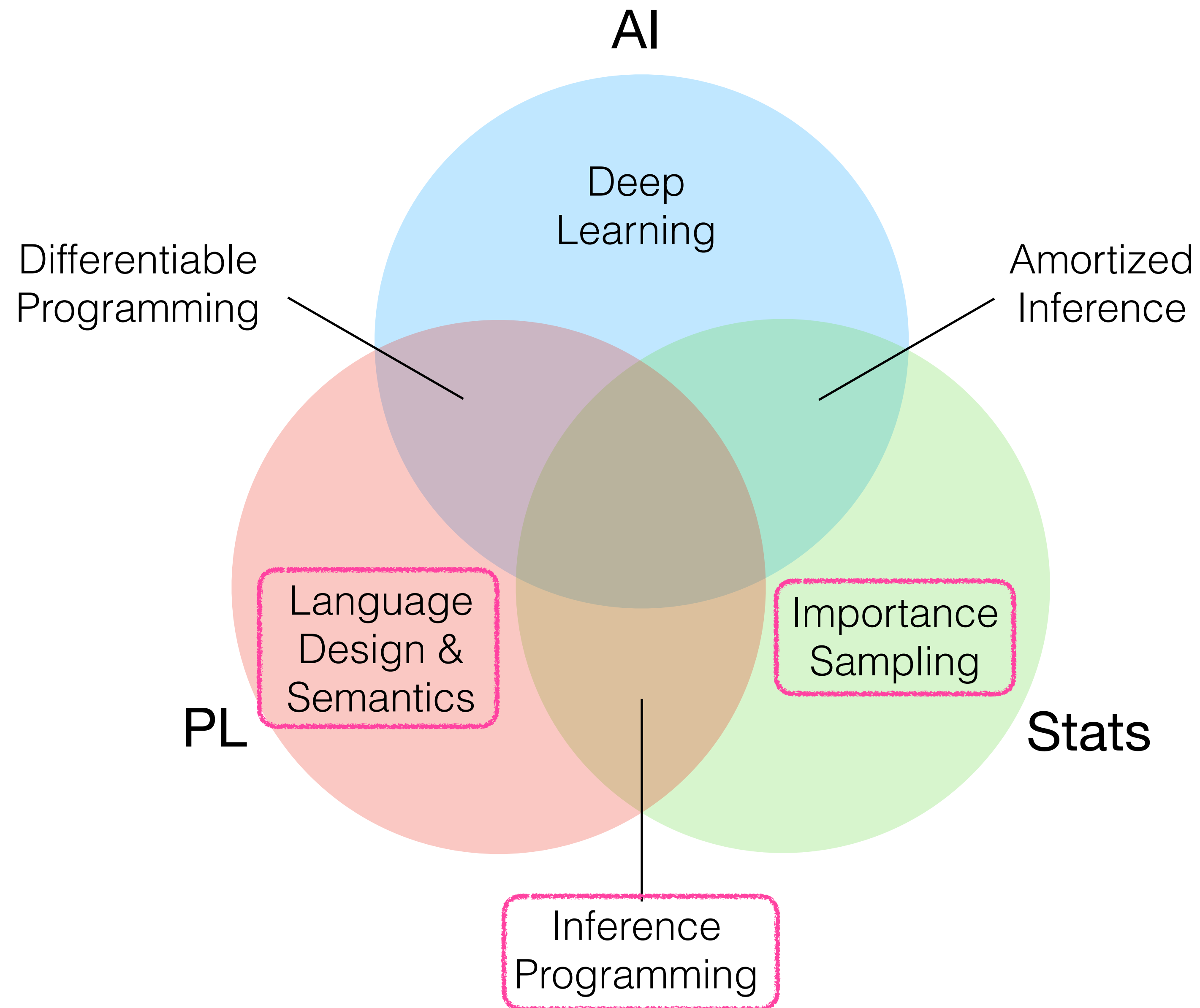
Inferred Positions

Reconstructions



- Completely unsupervised
- Computationally efficient (~5 updates, ~10 particles)

# Deep Probabilistic Programming



# Reasoning Compositionally About Inference

What (inference) DSL could define this sampler / variational method?

---

**Algorithm 1** Amortized Population Gibbs Sampling

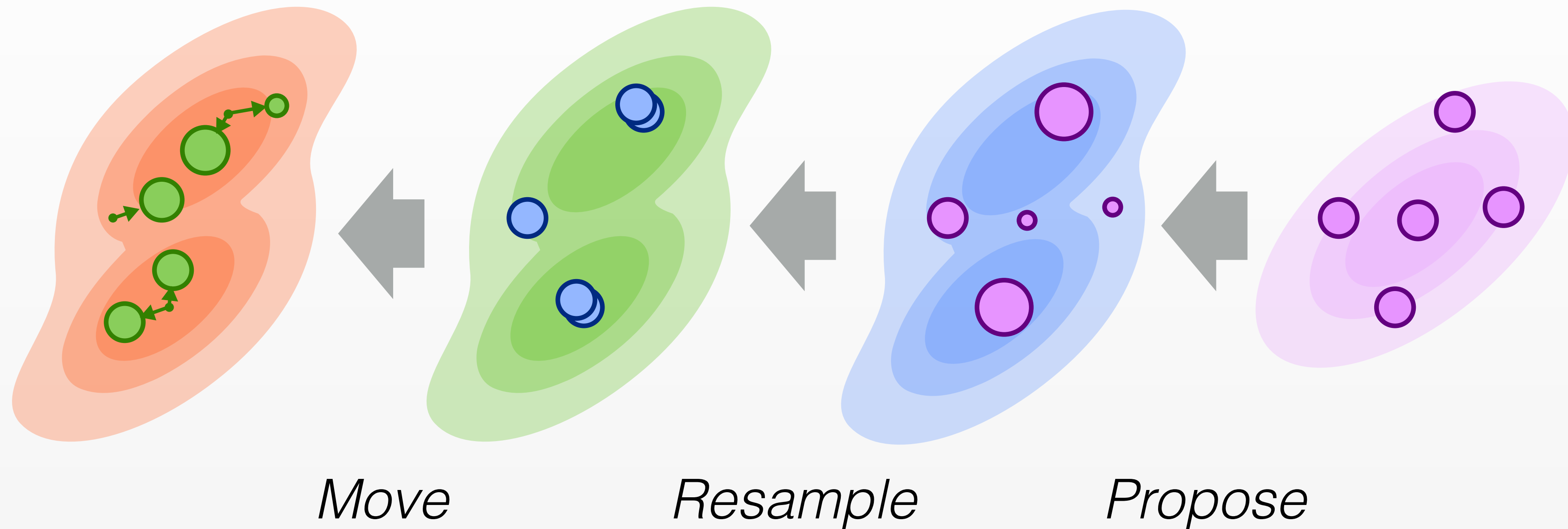
---

```
1: for  $n$  in  $1, \dots, N$  do
2:    $G_\phi = 0$ 
3:    $x^n \sim p^{\text{DATA}}(x)$ 
4:   for  $l$  in  $1, \dots, L$  do
5:      $z^{n,1,l} \sim q_\phi(z | x^n)$ 
6:      $w^{n,1,l} \leftarrow p_\theta(x^n, z^{n,1,l}) / q_\phi(z^{n,1,l})$ 
7:   for  $k$  in  $2, \dots, K$  do
8:      $\tilde{z}, \tilde{w} = z^{n,k-1}, w^{n,k-1}$ 
9:     for  $b$  in  $1, \dots, B$  do
10:       $\tilde{z}, \tilde{w} = \text{RESAMPLE}(\tilde{z}, \tilde{w})$ 
11:      for  $l$  in  $1, \dots, L$  do
12:         $\tilde{z}'_b{}^l \sim q_\phi(\cdot | x^n, \tilde{z}_{-b}^l)$ 
13:         $\tilde{w}^l = \frac{p_\theta(x^n, \tilde{z}'_b{}^l, \tilde{z}_{-b}^l) q_\phi(\tilde{z}'_b{}^l | x^n, \tilde{z}_{-b}^l)}{p_\theta(x^n, \tilde{z}_b^l, \tilde{z}_{-b}^l) q_\phi(\tilde{z}_b^l | x^n, \tilde{z}_{-b}^l)} \tilde{w}^l$ 
14:         $\tilde{z}_b^l = \tilde{z}'_b{}^l$ 
15:       $G_\phi = G_\phi + \sum_{l=1}^L \frac{\tilde{w}^l}{\sum_{l'} \tilde{w}^{l'}} \frac{d}{d\phi} \log q_\phi(\tilde{z}_b^l | x^n, \tilde{z}_{-b}^l)$ 
16:       $z^{n,k}, w^{n,k} = \tilde{z}, \tilde{w}$ 
17: return  $G_\phi, z, w$  ▷ Output: Grad
```

---

- APG is an example of a amortized SMC sampler
- Known building blocks, but not trivial to combine correctly
- Can we define compositional methods for importance sampling and gradient estimation?

# Combinators: A DSL for Inference

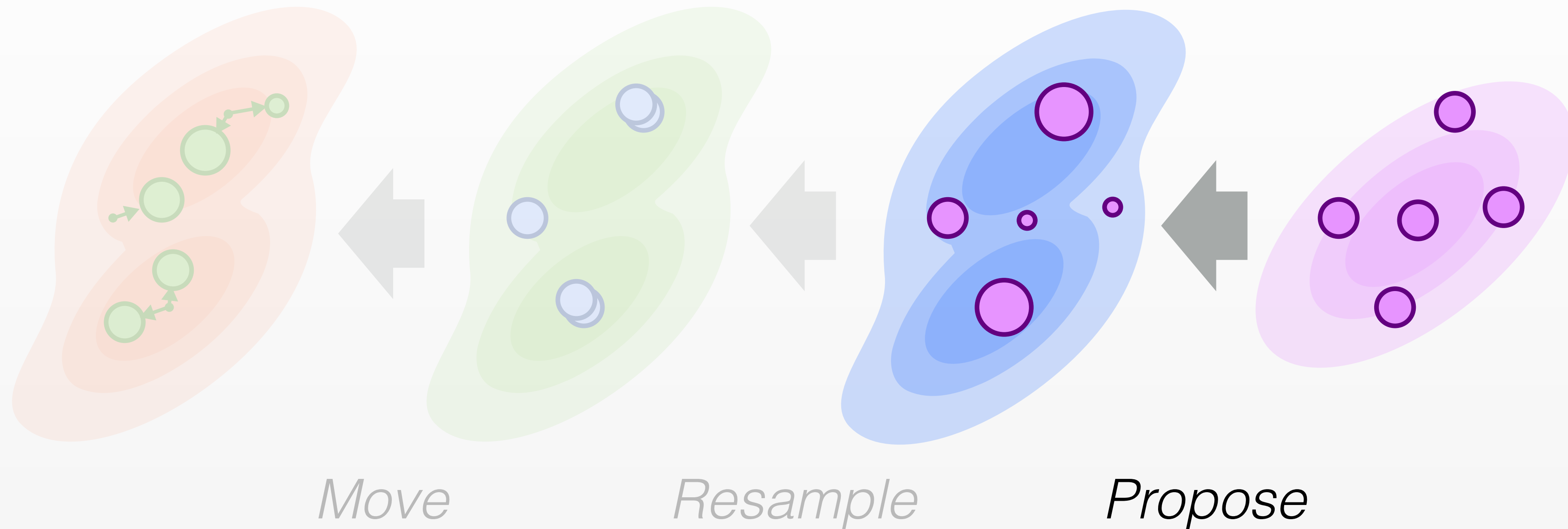


$f ::=$  A primitive program

$p ::= f \mid \text{extend}(p, f)$

$q ::= p \mid \text{resample}(q) \mid \text{compose}(q', q) \mid \text{propose}(p, q)$

# Combinators: A DSL for Inference



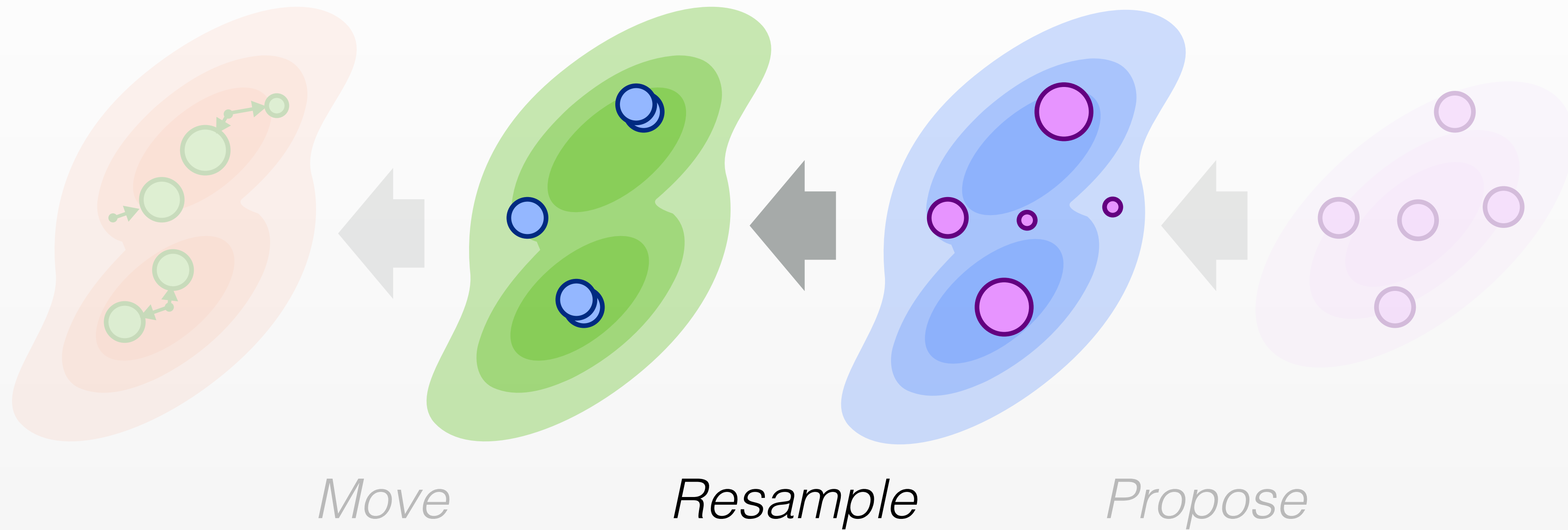
$f ::=$  A primitive program

$p ::= f \mid \text{extend}(p, f)$

$q ::= p \mid \text{resample}(q) \mid \text{compose}(q', q) \mid \text{propose}(p, q)$



# Combinators: A DSL for Inference

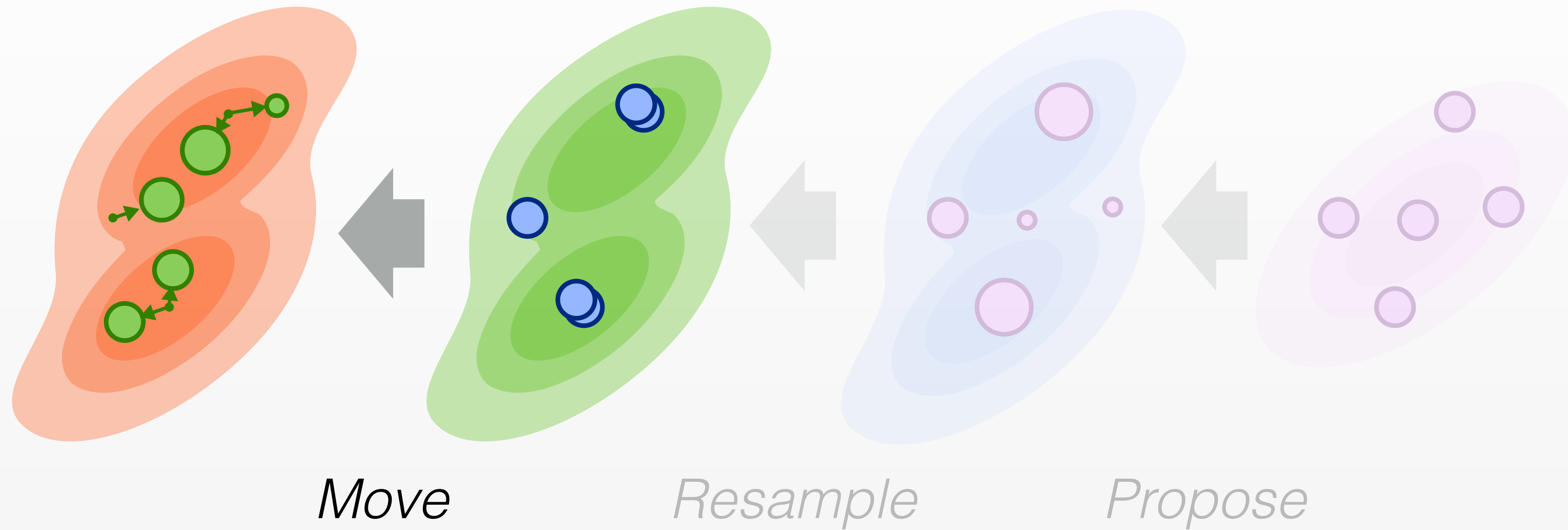


$f ::=$  A primitive program

$p ::= f \mid \text{extend}(p, f)$

$q ::= p \mid \text{resample}(q) \mid \text{compose}(q', q) \mid \text{propose}(p, q)$

# Combinators: A DSL for Inference

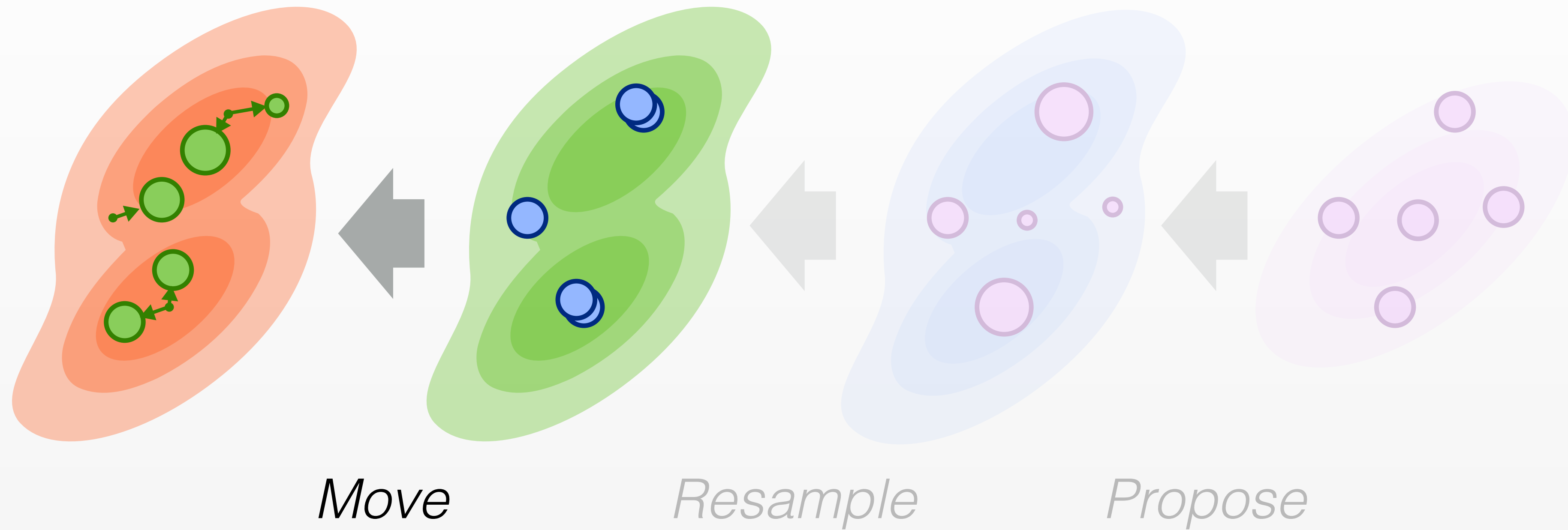


$f ::=$  A primitive program

$p ::= f \mid \text{extend}(p, f)$

$q ::= p \mid \text{resample}(q) \mid \text{compose}(q', q) \mid \text{propose}(p, q)$

# Combinators: A DSL for Inference



$f ::=$  A primitive program

$p ::= f \mid \text{extend}(p, f)$

$q ::= p \mid \text{resample}(q) \mid \text{compose}(q', q) \mid \text{propose}(p, q)$

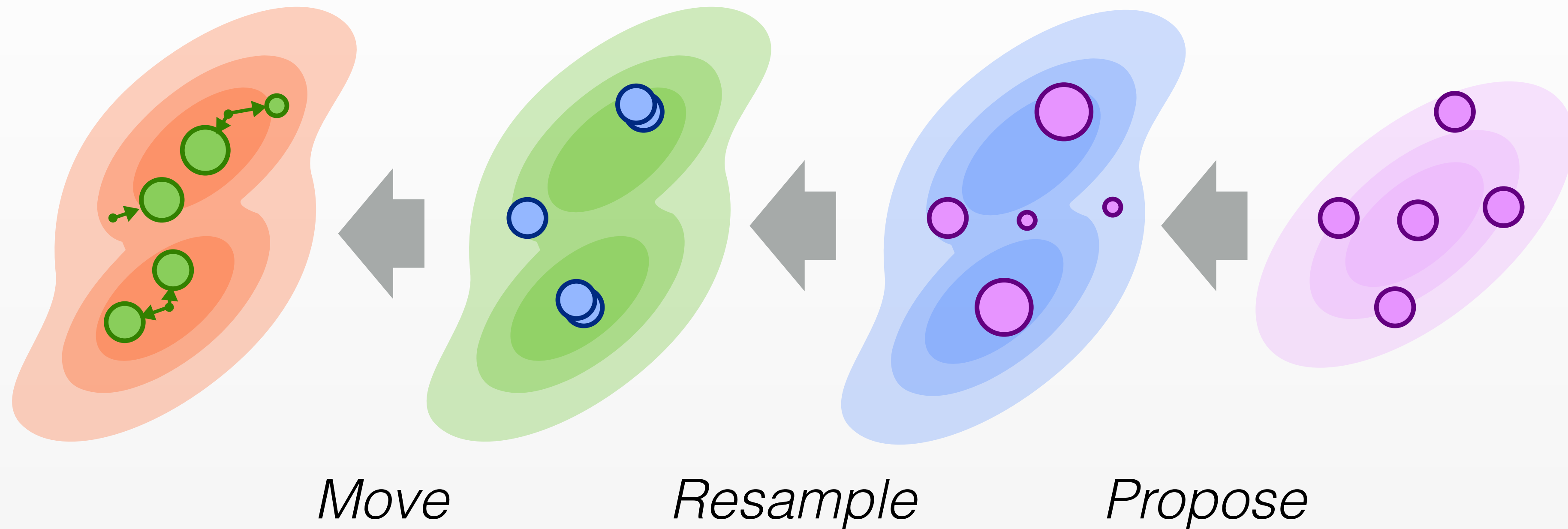
# Core Property: Proper Weighting

**Definition:** A pair  $w, \eta$  is properly weighted with respect to a density  $\pi(\eta)$  when, for all measurable  $h(\eta)$ ,

$$\mathbb{E}_{w, \eta \sim \Pi} [w h(\eta)] = \mathcal{C} \mathbb{E}_{\eta \sim \pi} [h(\eta)]$$



# Combinators: A DSL for Inference



**Semantics:** Composition preserves proper weighting

<https://github.com/probtorch/combinators>

(Pyro implementation coming this Summer)

# Example: Amortized Gibbs Samplers

```
def pop_gibbs(target, proposal, kernels, sweeps):  
    q = propose(partial(target, suffix=0),  
               partial(proposal, suffix=0))  
    for s in range(sweeps):  
        for k in kernels:  
            q = propose(  
                extend(partial(target, suffix=s+1),  
                      partial(k, suffix=s)),  
                compose(partial(k, suffix=s+1),  
                       resample(q, dim=0)))  
    return q
```

High-level algorithm description  
(transition operators, resampling)

---

**Algorithm 1** Amortized Population Gibbs Sampling

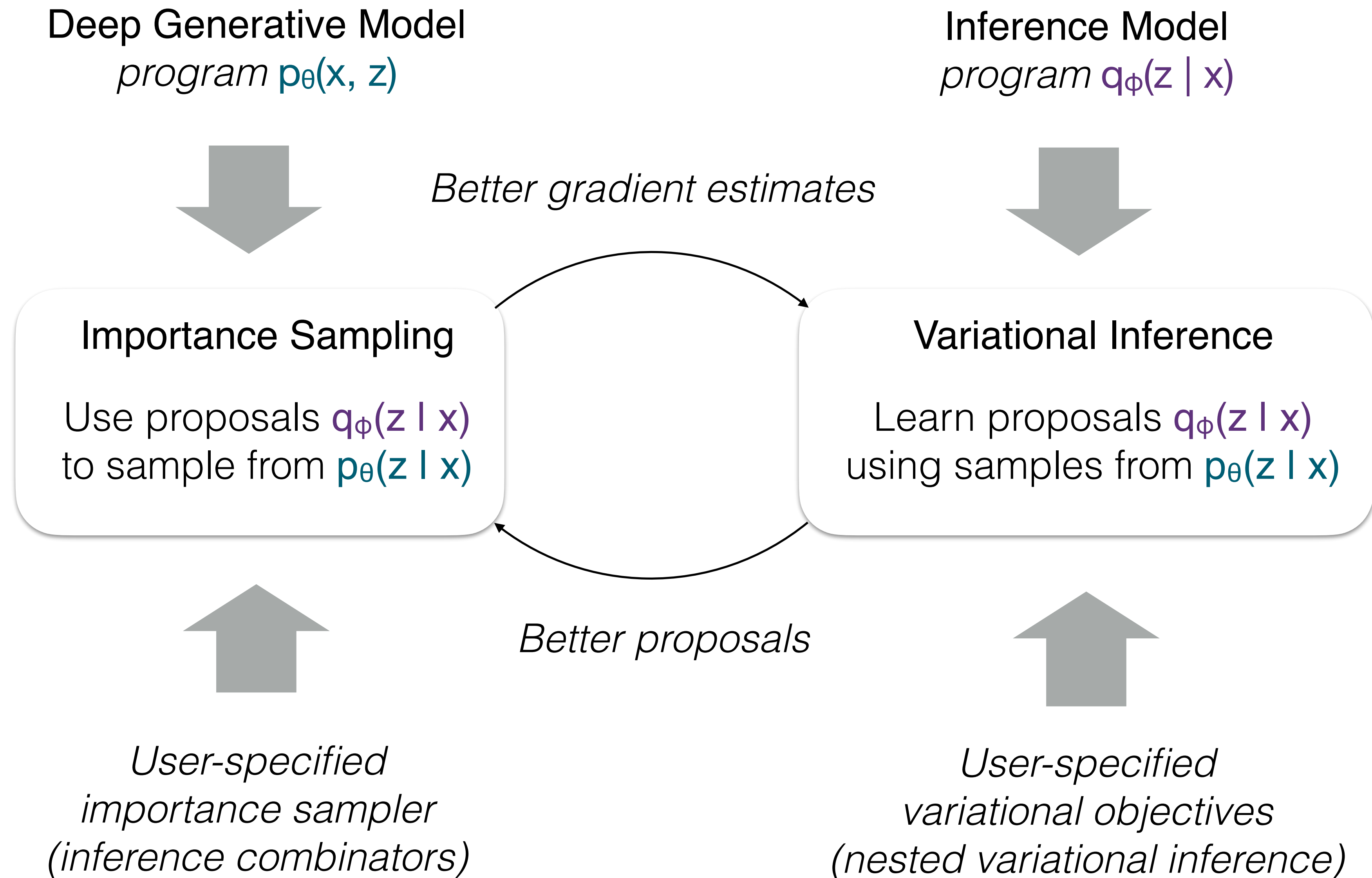
---

```
1: for  $n$  in  $1, \dots, N$  do  
2:    $G_\phi = 0$   
3:    $x^n \sim p^{\text{DATA}}(x)$   
4:   for  $l$  in  $1, \dots, L$  do  
5:      $z^{n,1,l} \sim q_\phi(z | x^n)$   
6:      $w^{n,1,l} \leftarrow p_\theta(x^n, z^{n,1,l}) / q_\phi(z^{n,1,l})$   
7:   for  $k$  in  $2, \dots, K$  do  
8:      $\tilde{z}, \tilde{w} = z^{n,k-1}, w^{n,k-1}$   
9:     for  $b$  in  $1, \dots, B$  do  
10:       $\tilde{z}, \tilde{w} = \text{RESAMPLE}(\tilde{z}, \tilde{w})$   
11:      for  $l$  in  $1, \dots, L$  do  
12:         $\tilde{z}_b^l \sim q_\phi(\cdot | x^n, \tilde{z}_{-b}^l)$   
13:         $\tilde{w}^l = \frac{p_\theta(x^n, \tilde{z}_b^l, \tilde{z}_{-b}^l) q_\phi(\tilde{z}_b^l | x^n, \tilde{z}_{-b}^l)}{p_\theta(x^n, \tilde{z}_b^l, \tilde{z}_{-b}^l) q_\phi(\tilde{z}_b^l | x^n, \tilde{z}_{-b}^l)} \tilde{w}^l$   
14:         $\tilde{z}_b^l = \tilde{z}_b^l$   
15:         $G_\phi = G_\phi + \sum_{l=1}^L \frac{\tilde{w}^l}{\sum_{l'} \tilde{w}^{l'}} \frac{d}{d\phi} \log q_\phi(\tilde{z}_b^l | x^n, \tilde{z}_{-b}^l)$   
16:       $z^{n,k}, w^{n,k} = \tilde{z}, \tilde{w}$   
17: return  $G_\phi, z, w$  ▷ Output: Grad
```

---

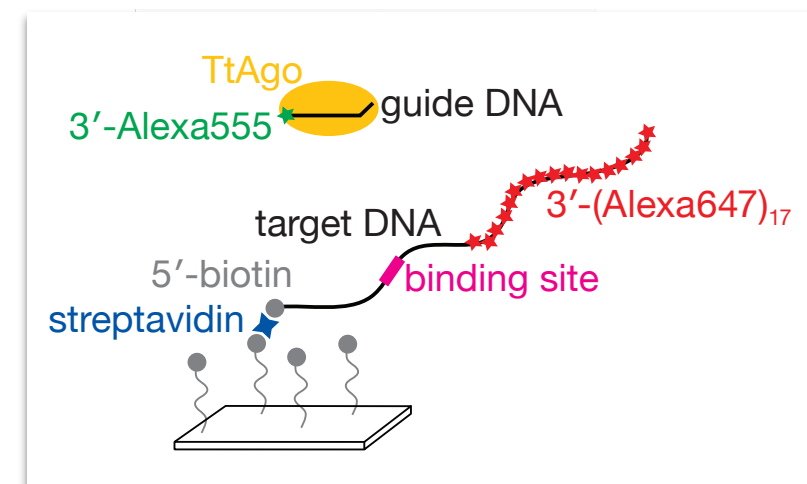
Low-level algorithm description  
(weight and gradient computations)

# Differentiable + Probabilistic + Inference Programming

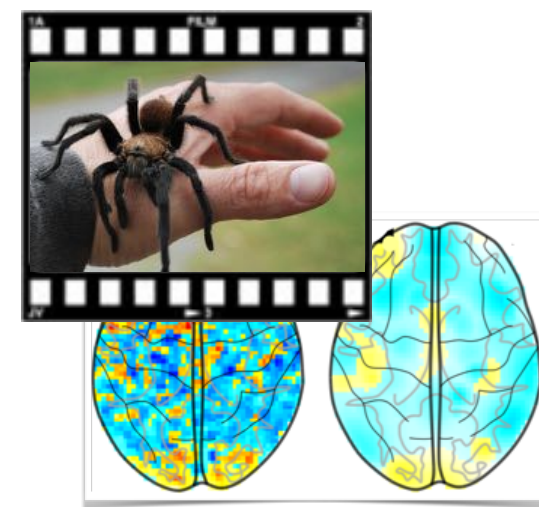


# The Next 700 Models in AI

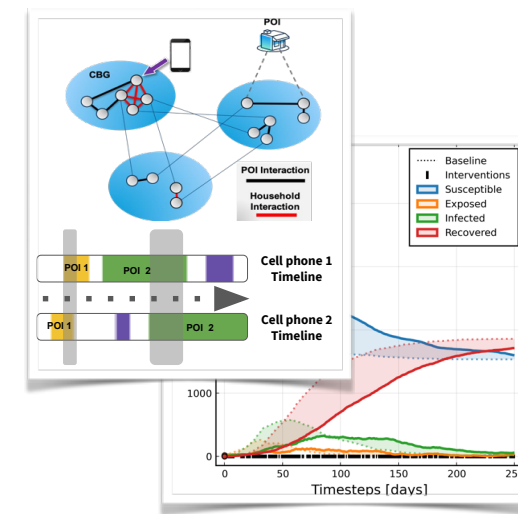
Biophysics



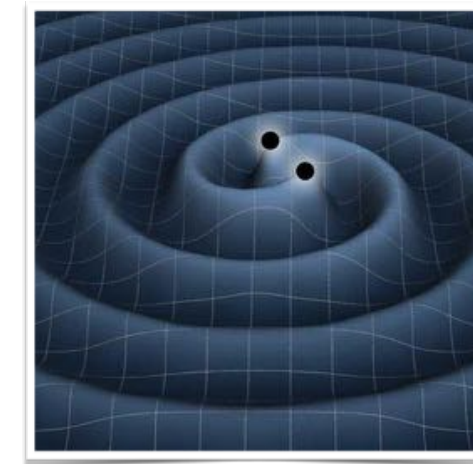
Neuroimaging



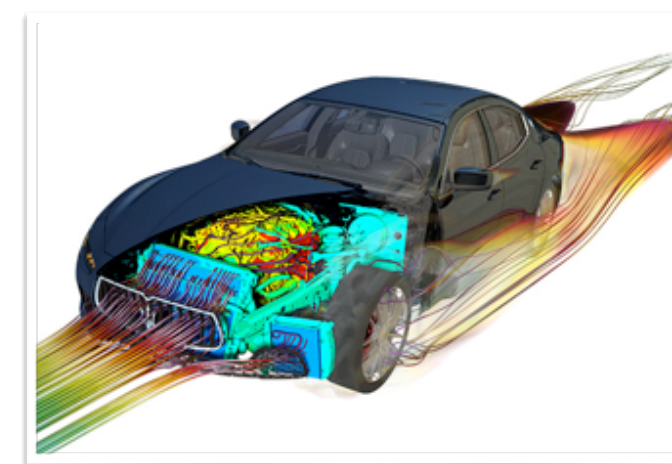
Epidemiology



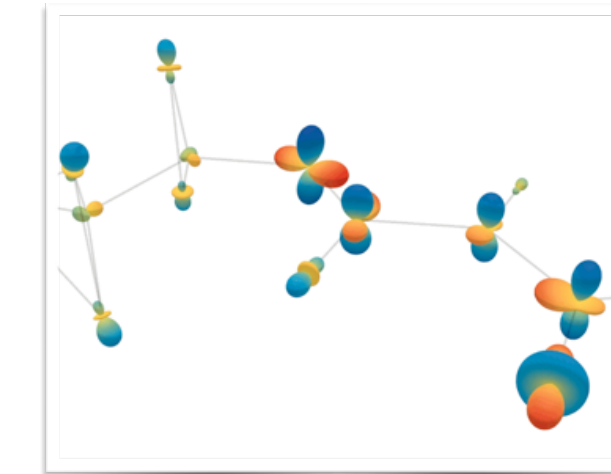
Astrophysics



Computational Fluid Dynamics



Molecular Design



Manufacturing



## Abstractions for Emerging Problems:

- Learning surrogate models
- Modeling search spaces
- Inferring differential equations



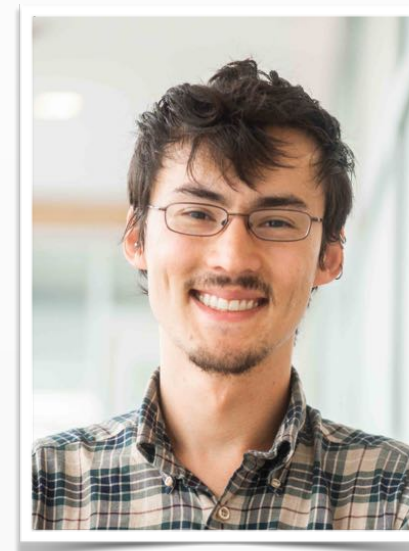
# Thank You!



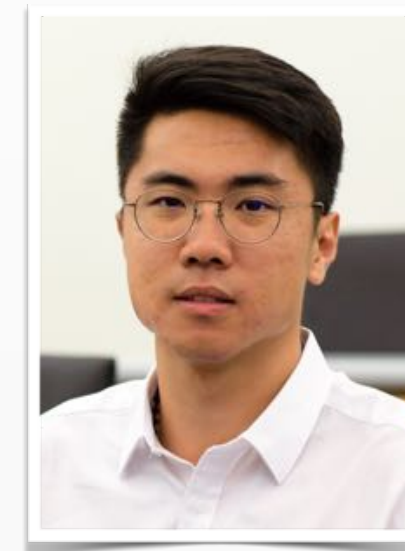
Heiko Zimmermann



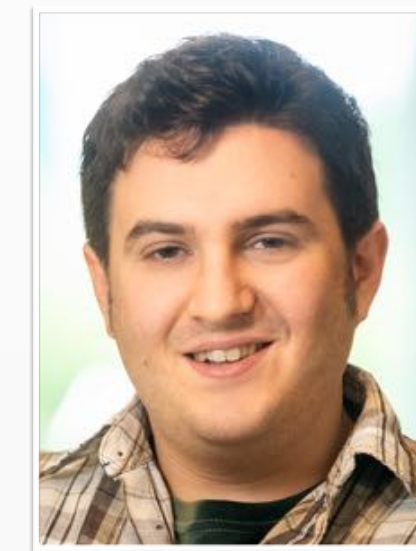
Babak Esmaeili



Sam Stites



Hao Wu



Eli Sennesh

## *Nested Variational Inference*

H. Zimmermann, H. Wu, Babak Esmaeili, J.-W. van de Meent  
NeurIPS 2021 [<https://arxiv.org/abs/2103.00668>]

## *Learning Proposals for Probabilistic Programs with Inference Combinators*

S. Stites\*, H. Zimmermann\*, H. Wu, E. Sennesh, J.-W. van de Meent  
UAI 2021 [<https://arxiv.org/abs/2103.00668>]

## *Amortized Population Gibbs Samplers with Neural Sufficient Statistics*

H. Wu, H. Zimmermann, E. Sennesh, Tuan Anh Le, J.-W. van de Meent  
ICML 2020 [[https://proceedings.icml.cc/static/paper\\_files/icml/2020/5881-Paper.pdf](https://proceedings.icml.cc/static/paper_files/icml/2020/5881-Paper.pdf)]

## *An Introduction to Probabilistic Programming*

J.-W. van de Meent, B. Paige, H. Yang, F. Wood  
ArXiv 2018 [<https://arxiv.org/abs/1809.10756>]