

DiCE: The Infinitely Differentiable Monte-Carlo Estimator Jakob Foerster¹, Gregory Farquhar¹*, Maruan Al-Shedivat²*, Tim Rocktäschel¹, Eric Xing², Shimon Whiteson¹

Motivation

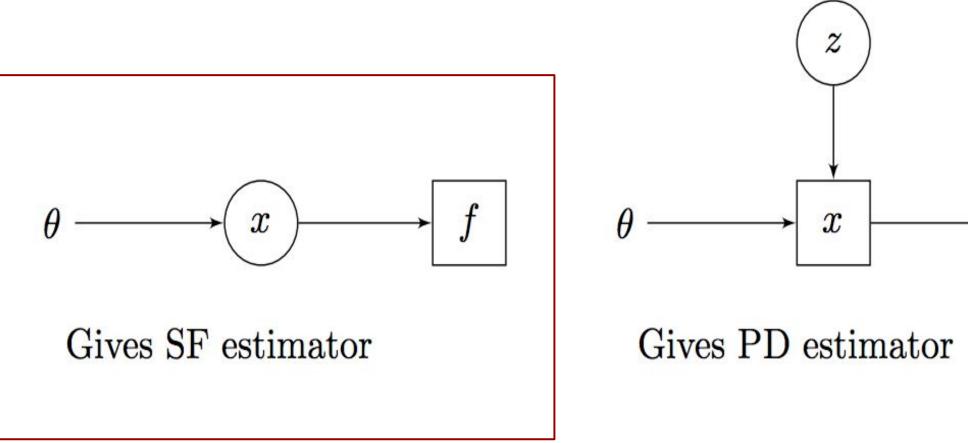
- The score function is commonly used in order to estimate gradients in Reinforcement Learning problems.
- This can be calculated via the Surrogate Loss approach.
- For order derivatives there is currently no satisfactory solution.
- Examples include Meta-Learning for RL and 2nd order optimisation.

Setting and background

• Stochastic Computation Graphs provide a framework for estimating gradients using *Surrogate Losses* (John Schulman et al 2015):

Input node Deterministic node

Stochastic node



• Surrogate Loss Approach:

$$SL(\Theta, S) := \sum_{w \in S} \log p(w \mid \text{DEPS}_w) \hat{Q}_w + \sum_{c \in C} c(\text{DEPS}_c)$$

Differentiation
$$In the treated independent of the treated indepndent of the tre$$

• Gradient Estimator:

$$\frac{\partial}{\partial \theta} \mathbb{E}\left[\sum_{c \in \mathcal{C}} c\right] = \mathbb{E}\left[\sum_{\substack{w \in \mathcal{S}, \\ \theta \prec^{D} w}} \left(\frac{\partial}{\partial \theta} \log p(w \mid \mathsf{DEPS}_w)\right) \hat{Q}_w + \sum_{\substack{c \in \mathcal{C} \\ \theta \prec^{D} c}} \frac{\partial}{\partial \theta} c(\mathsf{DEPS}_c)\right]\right]$$

- Surrogate Loss Approach cuts dependency of objective on parameters (^)
- This will cause errors in higher-order derivatives
- Compare with exact differentiation:

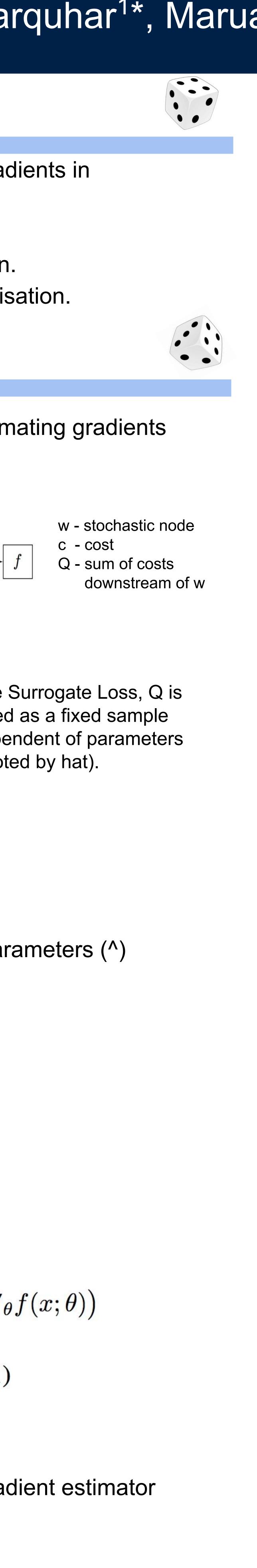
$$\begin{aligned} \nabla_{\theta} \mathcal{L} &= \nabla_{\theta} \mathbb{E}_{x} \left[f(x;\theta) \right] \\ &= \nabla_{\theta} \sum_{x} p(x;\theta) f(x;\theta) \\ &= \sum_{x} \nabla_{\theta} \left(p(x;\theta) f(x;\theta) \right) \\ &= \sum_{x} \left(f(x;\theta) \nabla_{\theta} p(x;\theta) + p(x;\theta) \nabla_{\theta} f(x;\theta) \right) \\ &= \sum_{x} \left(f(x;\theta) p(x;\theta) \nabla_{\theta} \log(p(x;\theta)) + p(x;\theta) \nabla_{\theta} \right) \\ &= \mathbb{E}_{x} \left[f(x;\theta) \nabla_{\theta} \log(p(x;\theta)) + \nabla_{\theta} f(x;\theta) \right] \quad (3.1) \\ &= \mathbb{E}_{x} \left[g(x;\theta) \right] \end{aligned}$$

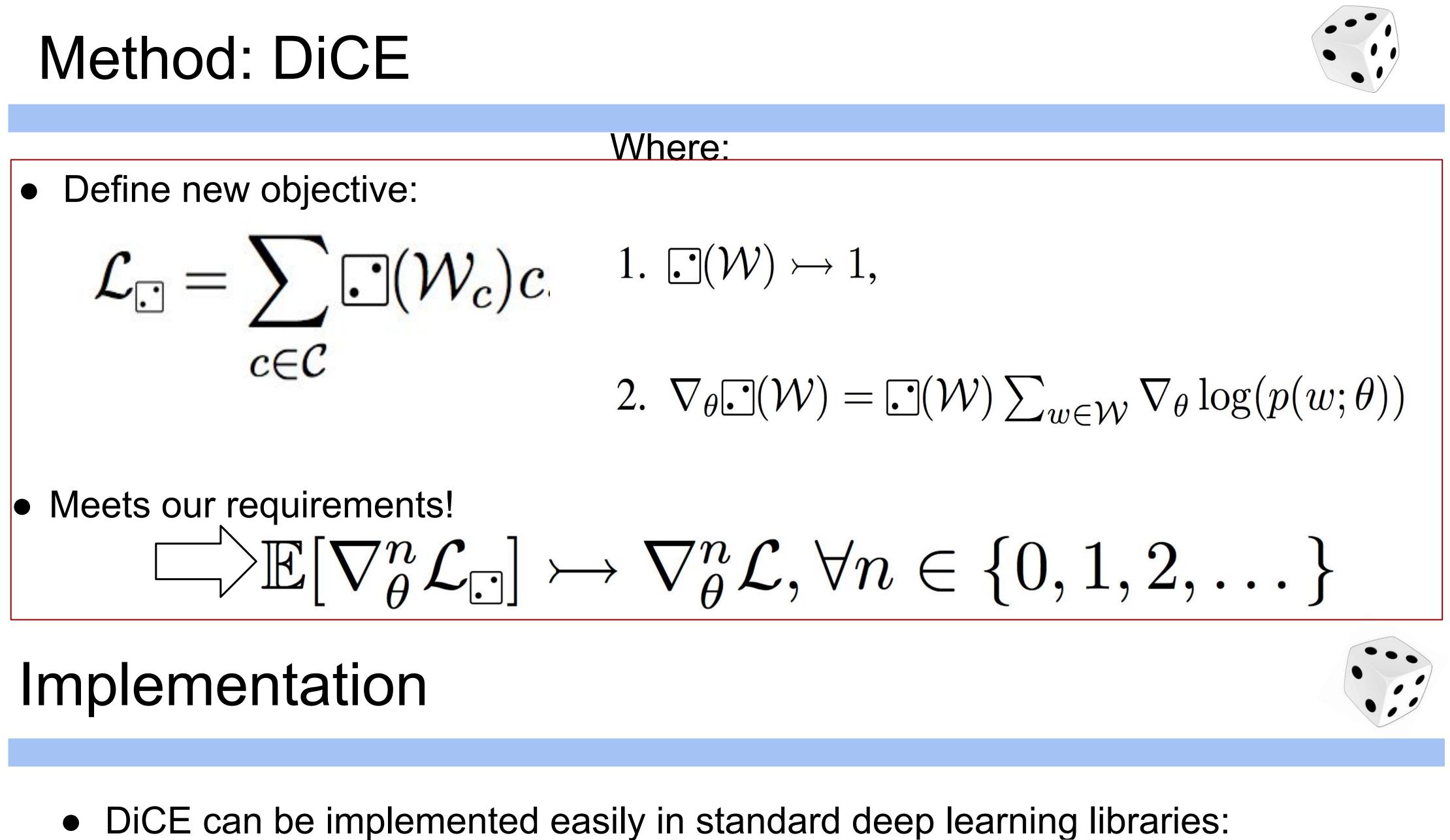
• Dependencies of $f(x;\theta)$ on θ should be maintained in the gradient estimator

So, are we done? Not quite...:

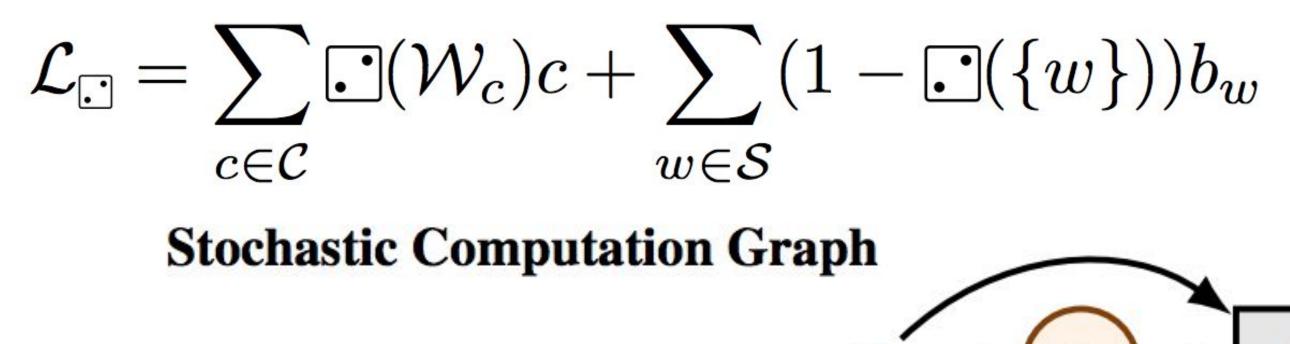
- How do we compute these gradients efficiently?
- Ideally would like to have an objective that can be differentiated multiple times.
- The gradient of the estimator should produce estimate of the gradient

¹University of Oxford, ²Carnegie Mellon University

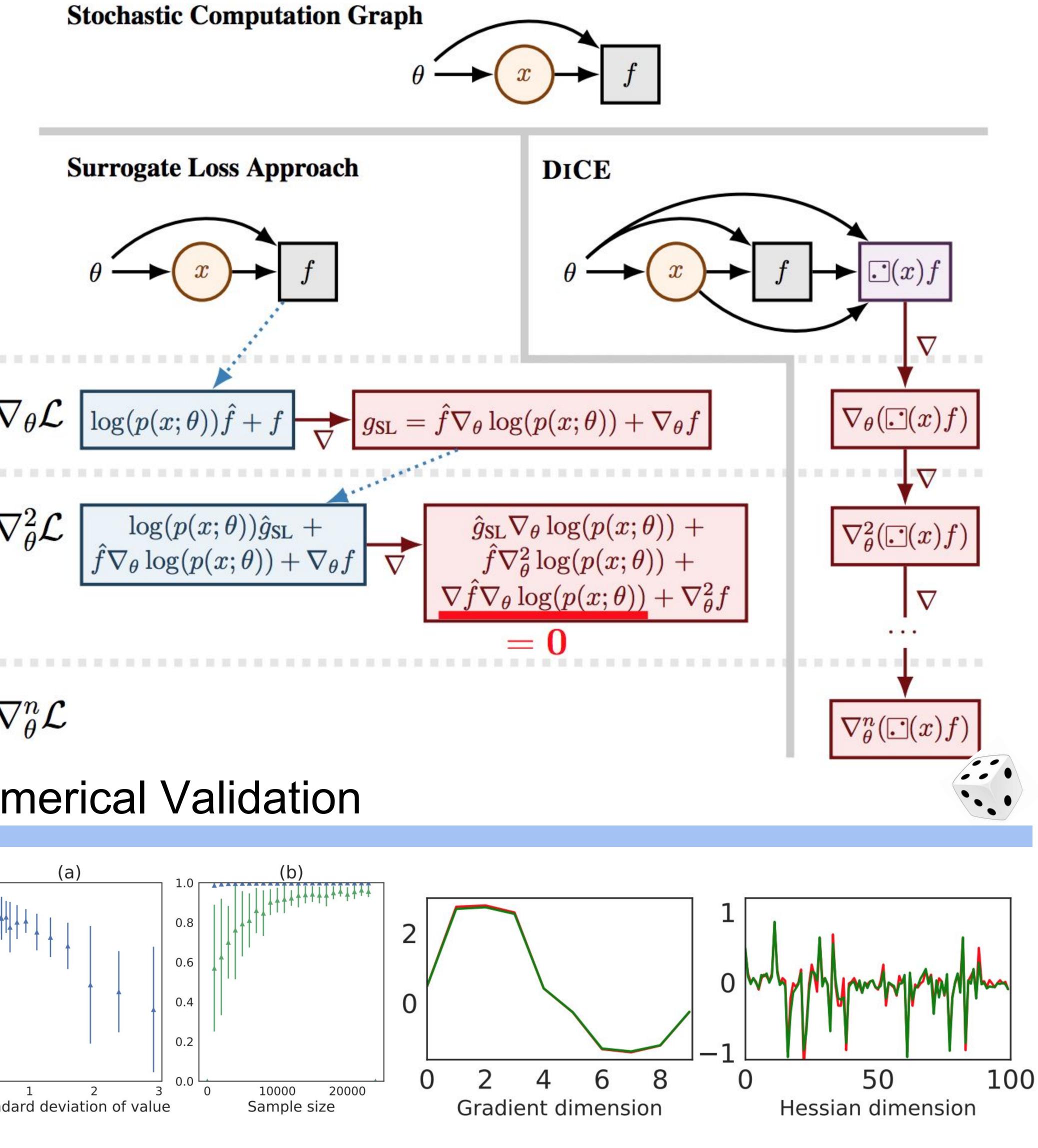


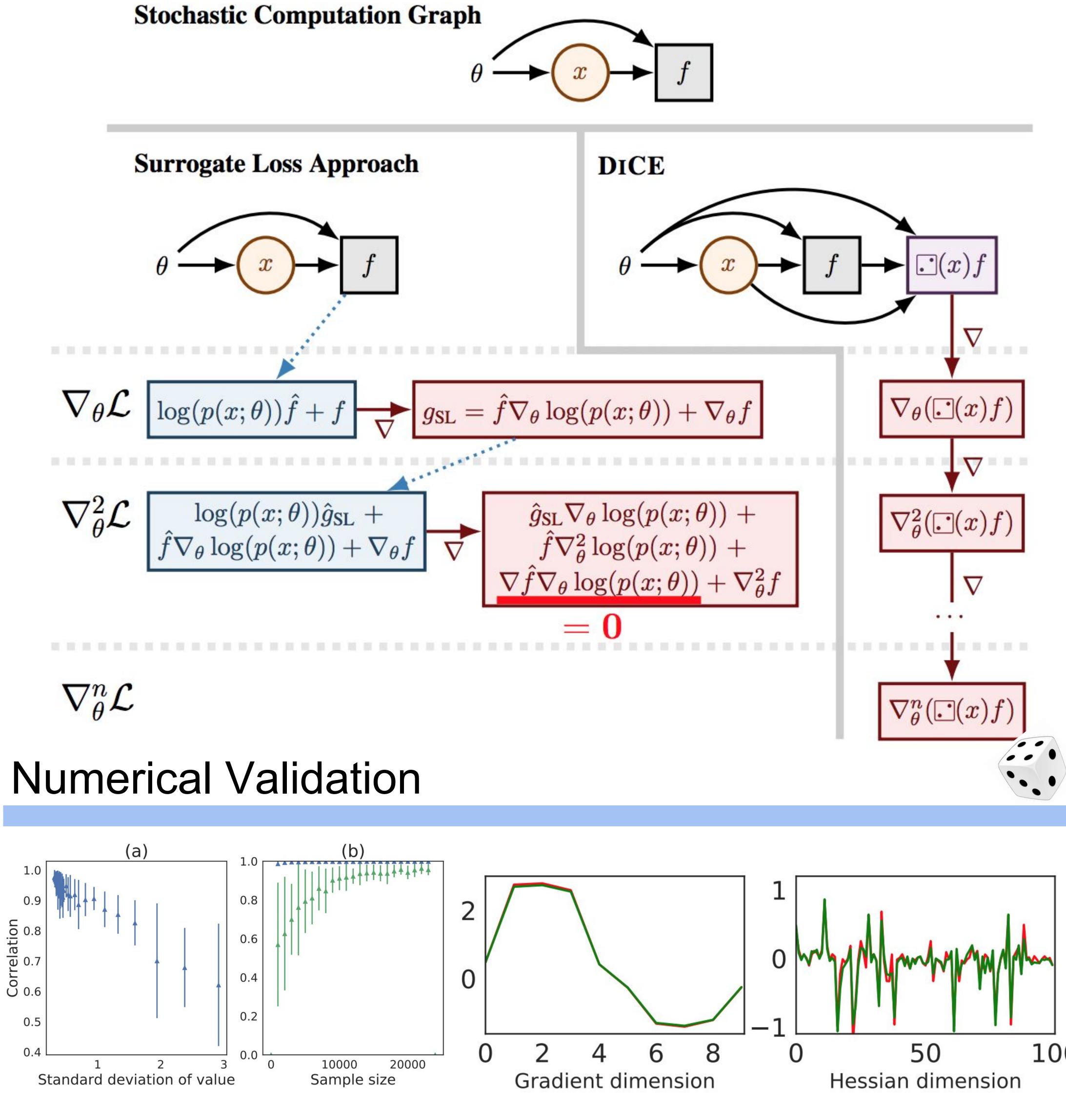


- Can include baseline naturally



Surrogate Loss Approach $\nabla_{\theta} \mathcal{L}$

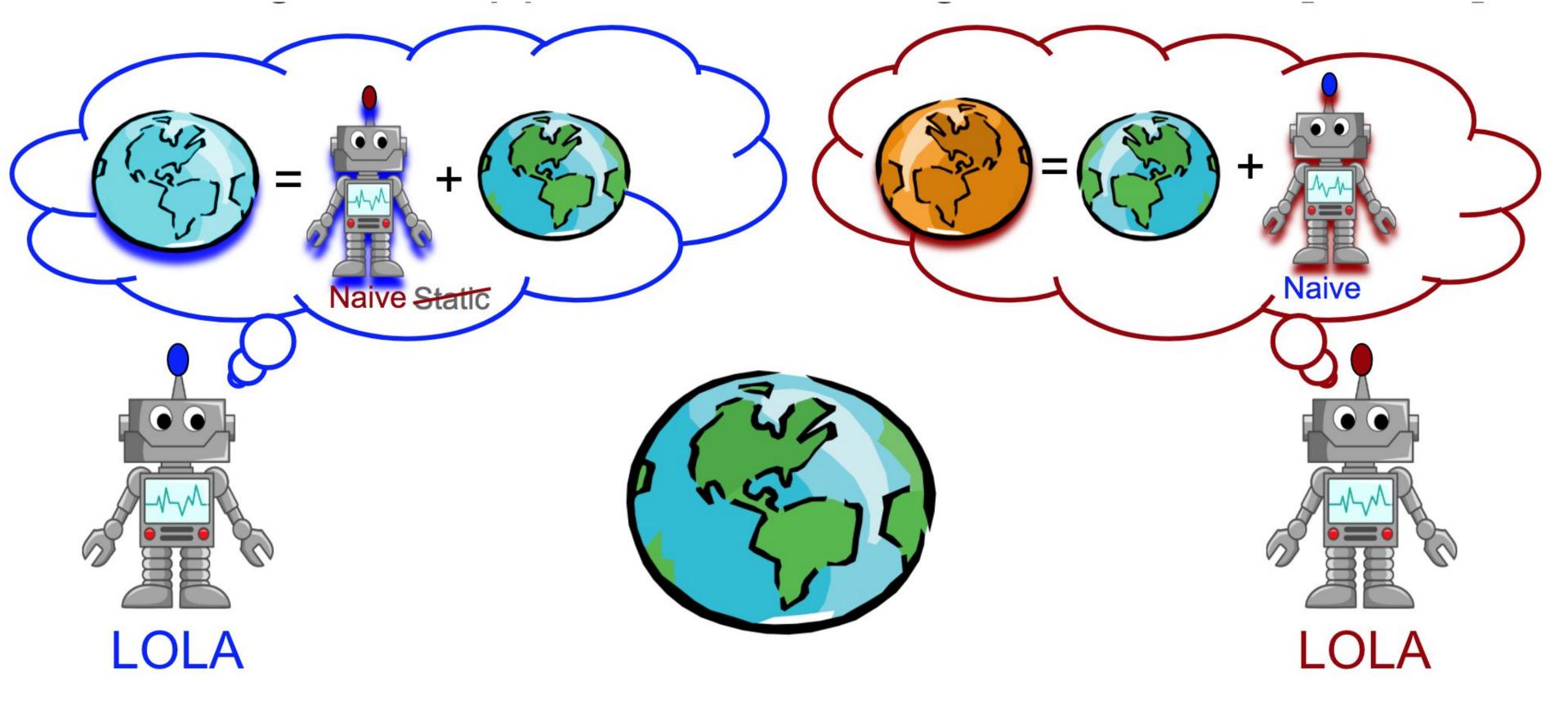




 $\Box(\mathcal{W}) = \exp\left(\tau - \bot(\tau)\right), \quad \tau = \sum \log(p(w;\theta)), \quad \bot(\tau): \mathsf{stop_gradient}$ $w {\in} \mathcal{W}$

Application Example: LOLA-DiCE

- When multiple agents are learning the learning step of each of the agents depends on the policy of all other agents.
- "Learning with Opponent-Learning Awareness" (Foerster et al 2017) explicitly accounts for this dependency and differentiates through the learning step:

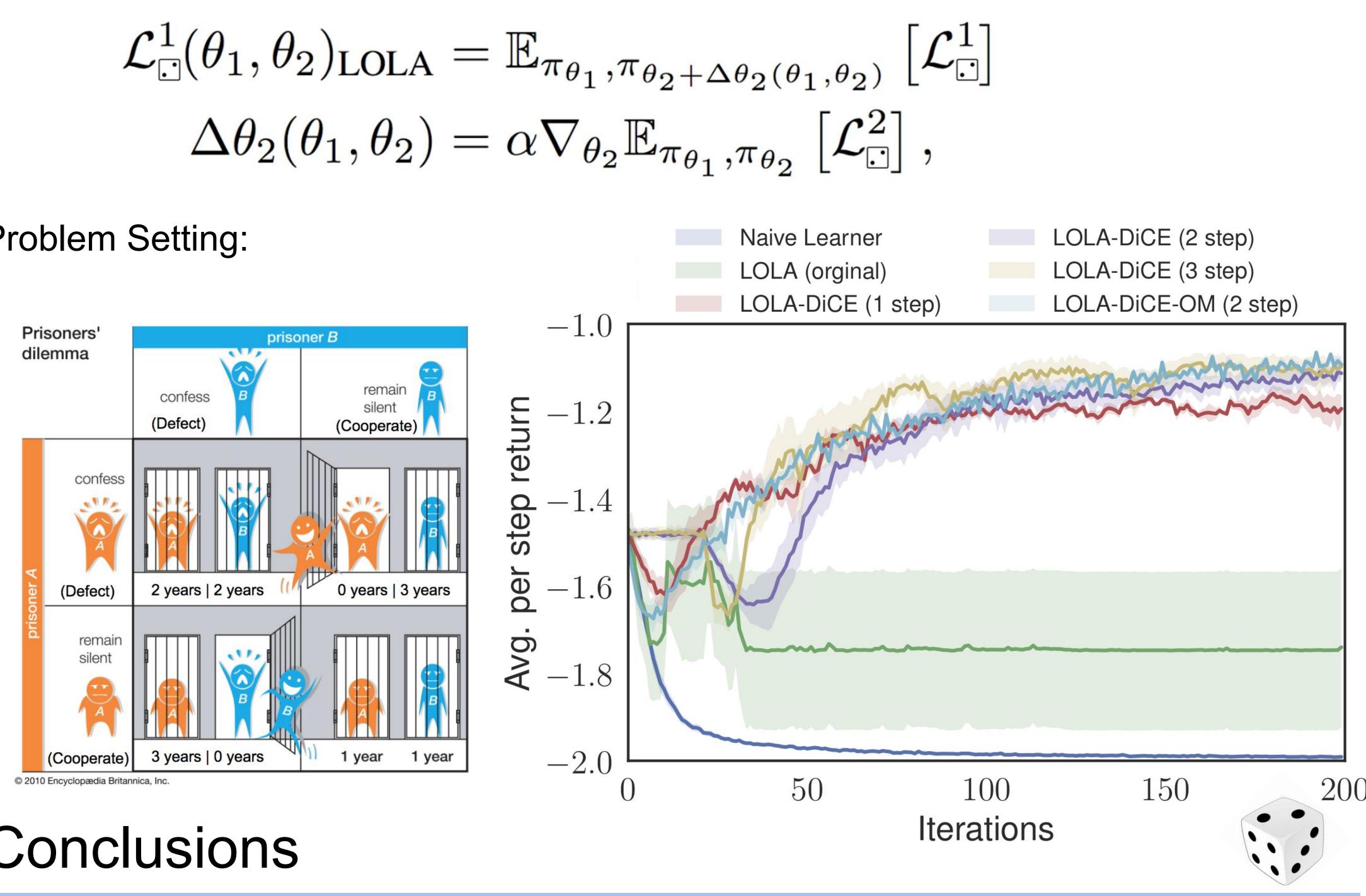


• The original paper relied on a Taylor expansion, second order gradients and a large batch size (4000). LOLA can be implemented exactly using DiCE:

$$\mathcal{L}^1_{\odot}(heta_1, heta_2)$$

 $\Delta heta_2(heta_2)$

Problem Setting:



Conclusions

- DiCE offers a new approach to estimating gradients in stochastic computation graphs • Arbitrary order gradients of the DiCE objective correspond to gradient estimators • DiCE can be implemented easily in standard learning frameworks!

Acknowledgements and References

Carnegie Vellon University



- Thanks to: Oxford-Google DeepMind Graduate Scholarship, the UK EPSRC CDT in Autonomous Intelligent Machines and Systems
- [1] "Gradient Estimation using Stochastic Computation Graphs" (John Schulman, Nicolas Heess, Theophane Weber, Pieter Abbeel)
- [2] "Learning with Opponent-Learning Awareness" (Jakob N. Foerster, Richard Y. Chen, Maruan Al-Shedivat, Shimon Whiteson, Pieter Abbeel, Igor Mordatch)