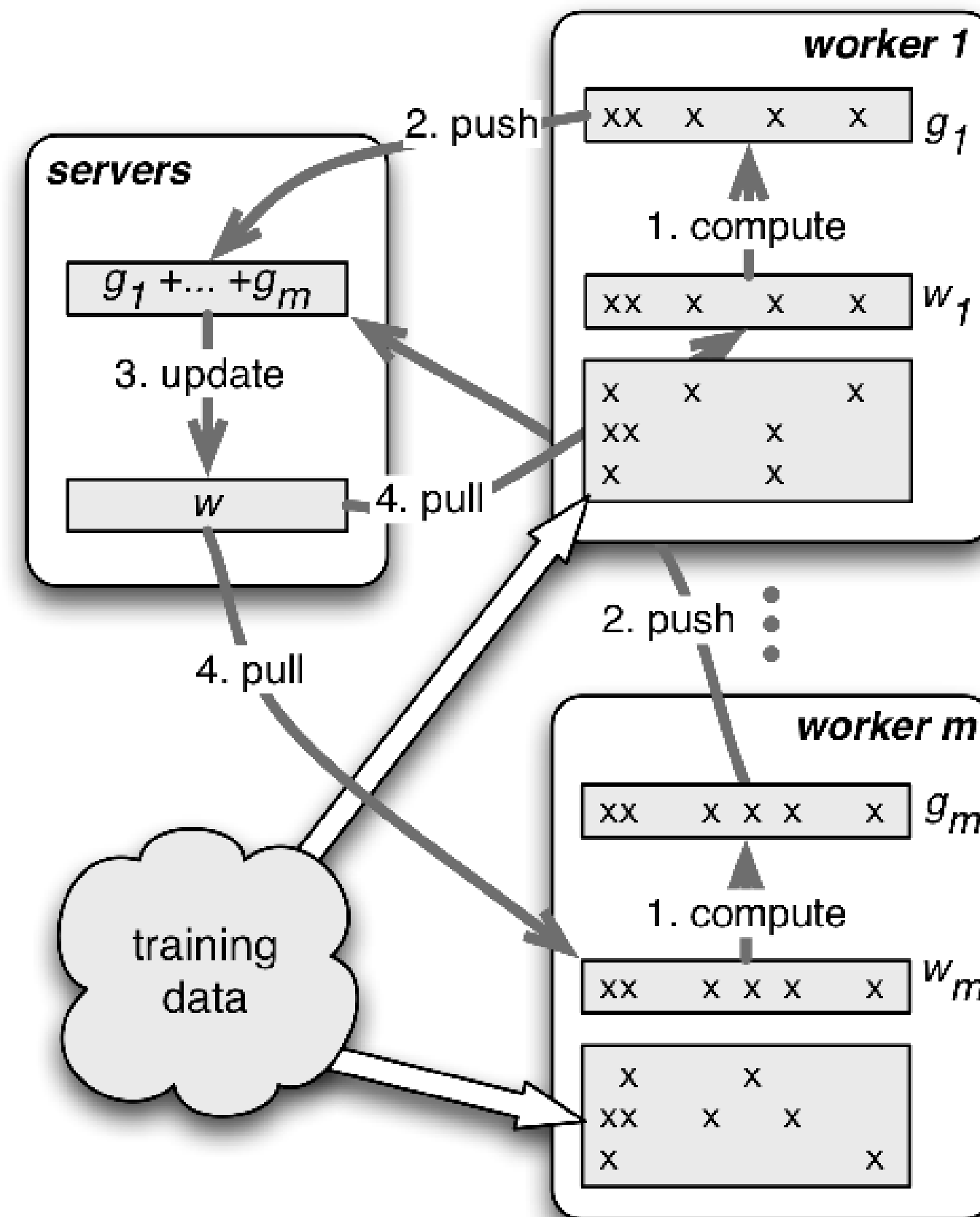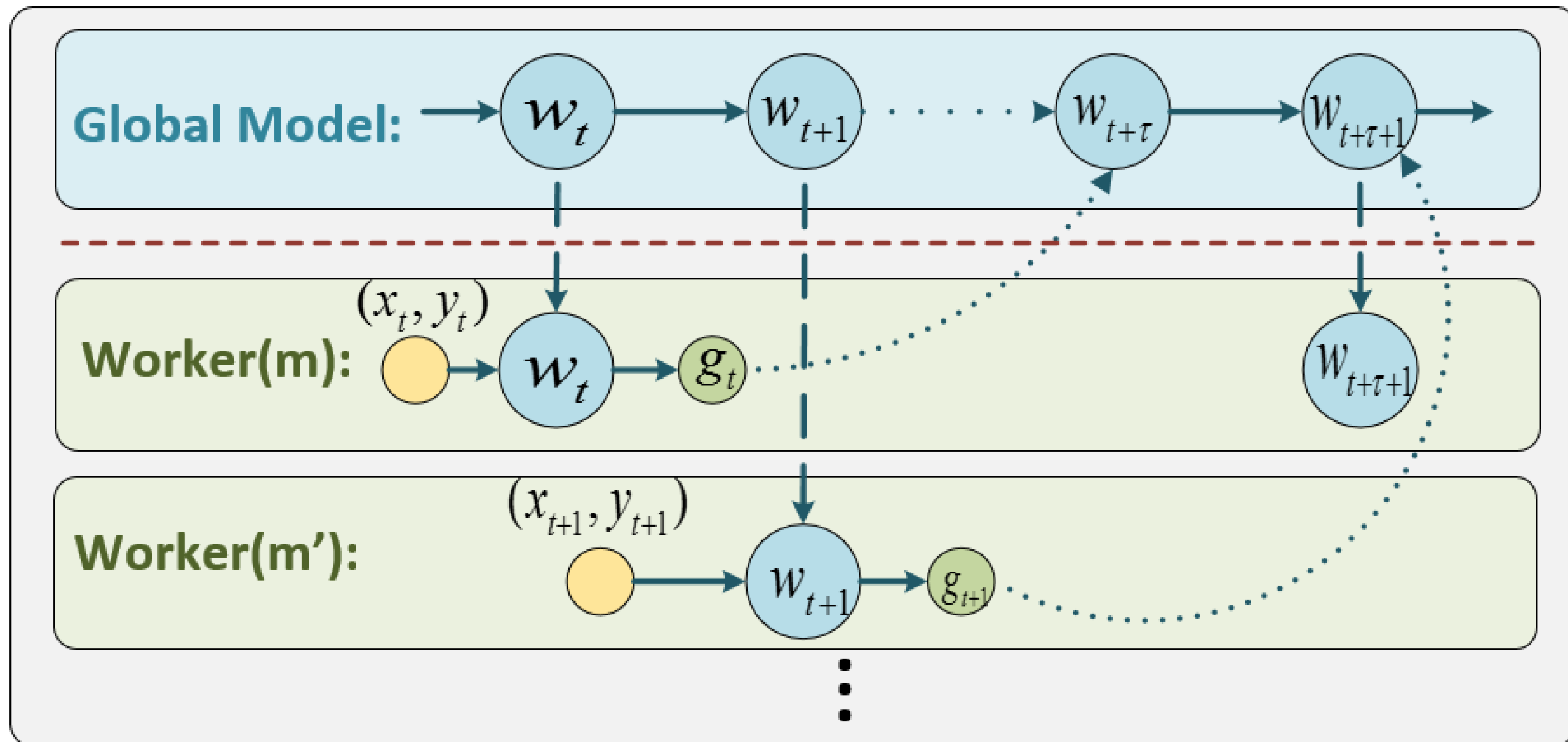# Bayesian Distributed Stochastic Gradient Descent

$$w = w - \eta \nabla Q(w) = w - \eta \sum_{i=1}^{\infty} \nabla Q_i(w)/n$$

# Parameter Server
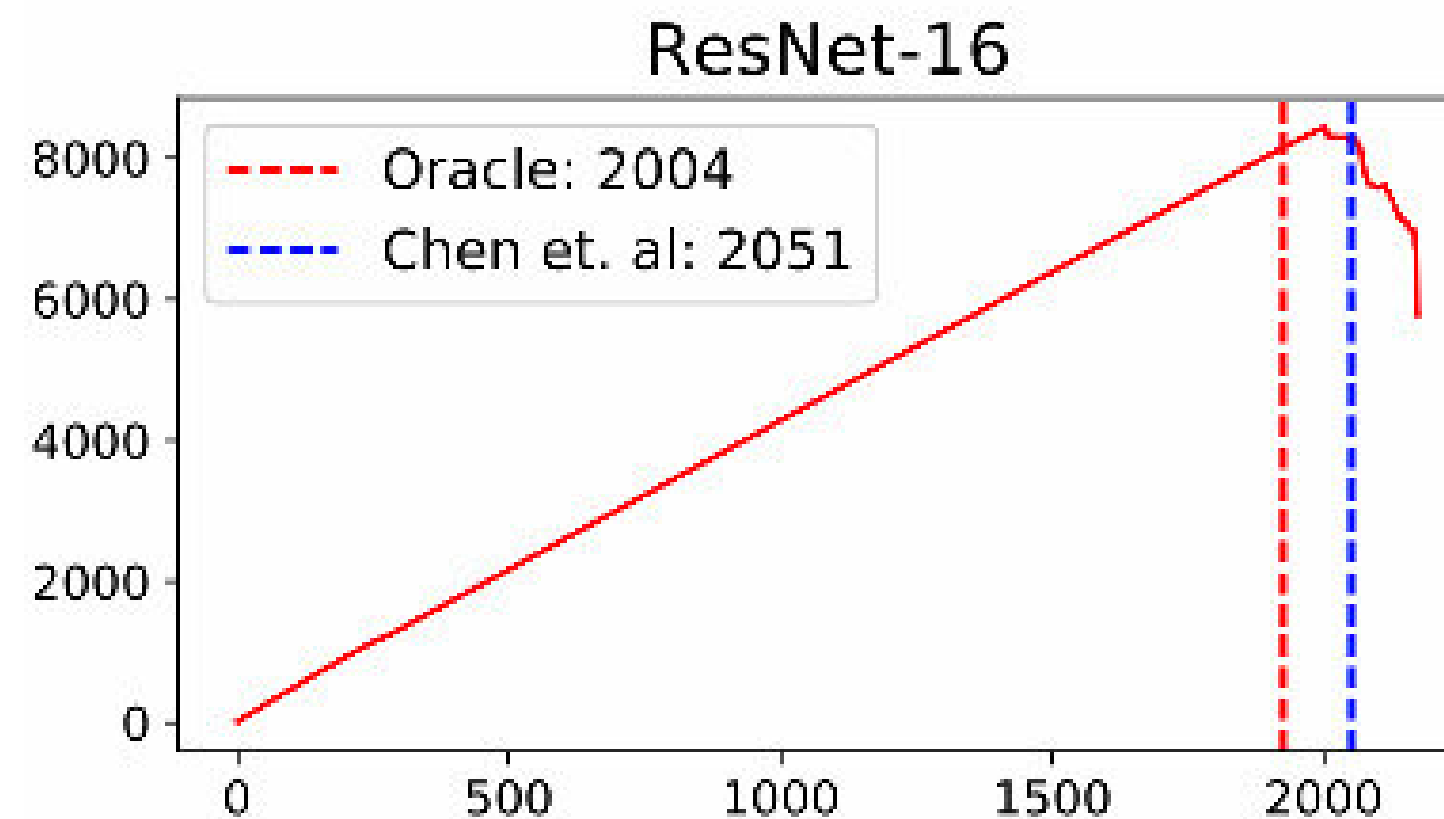
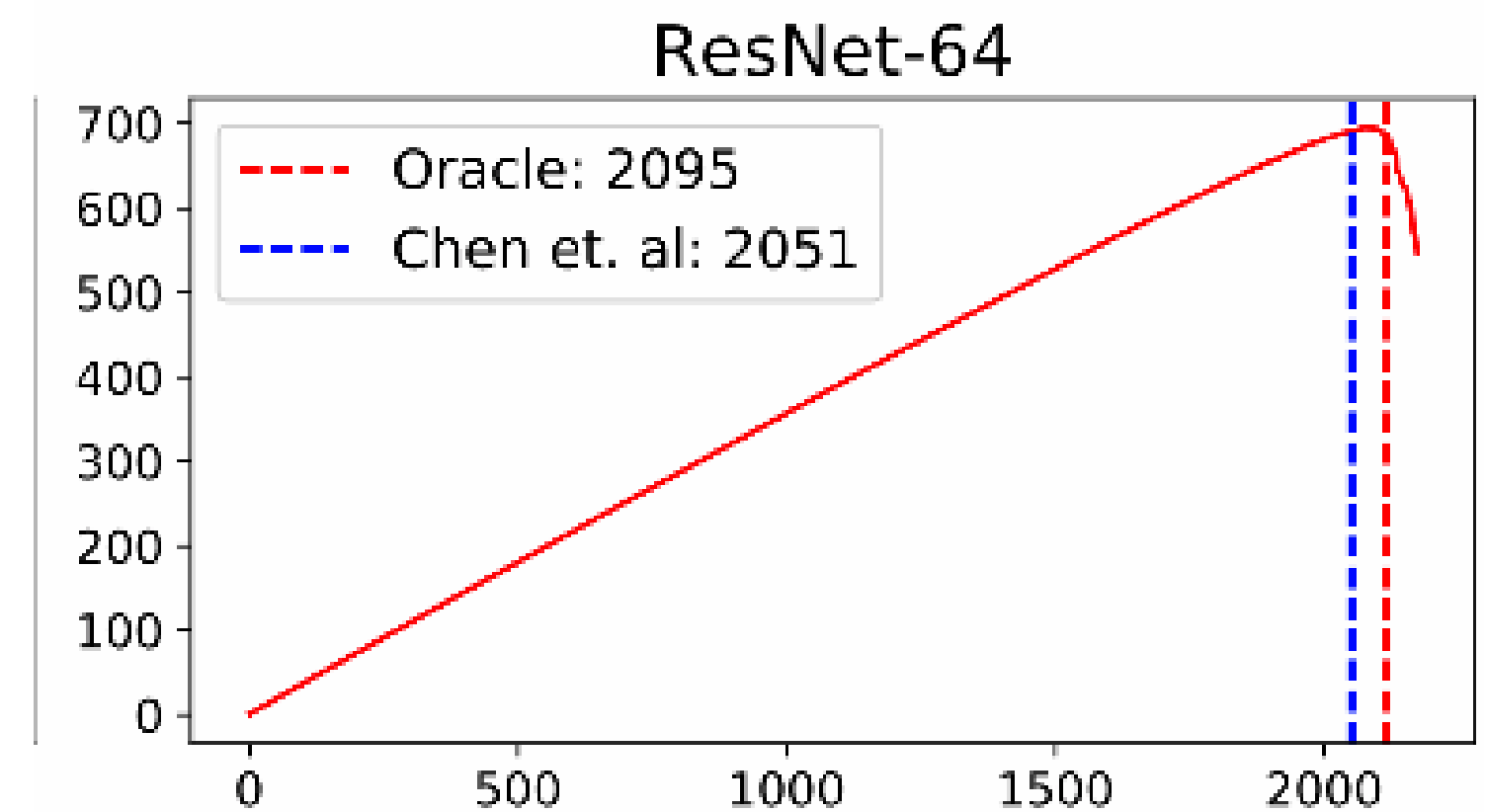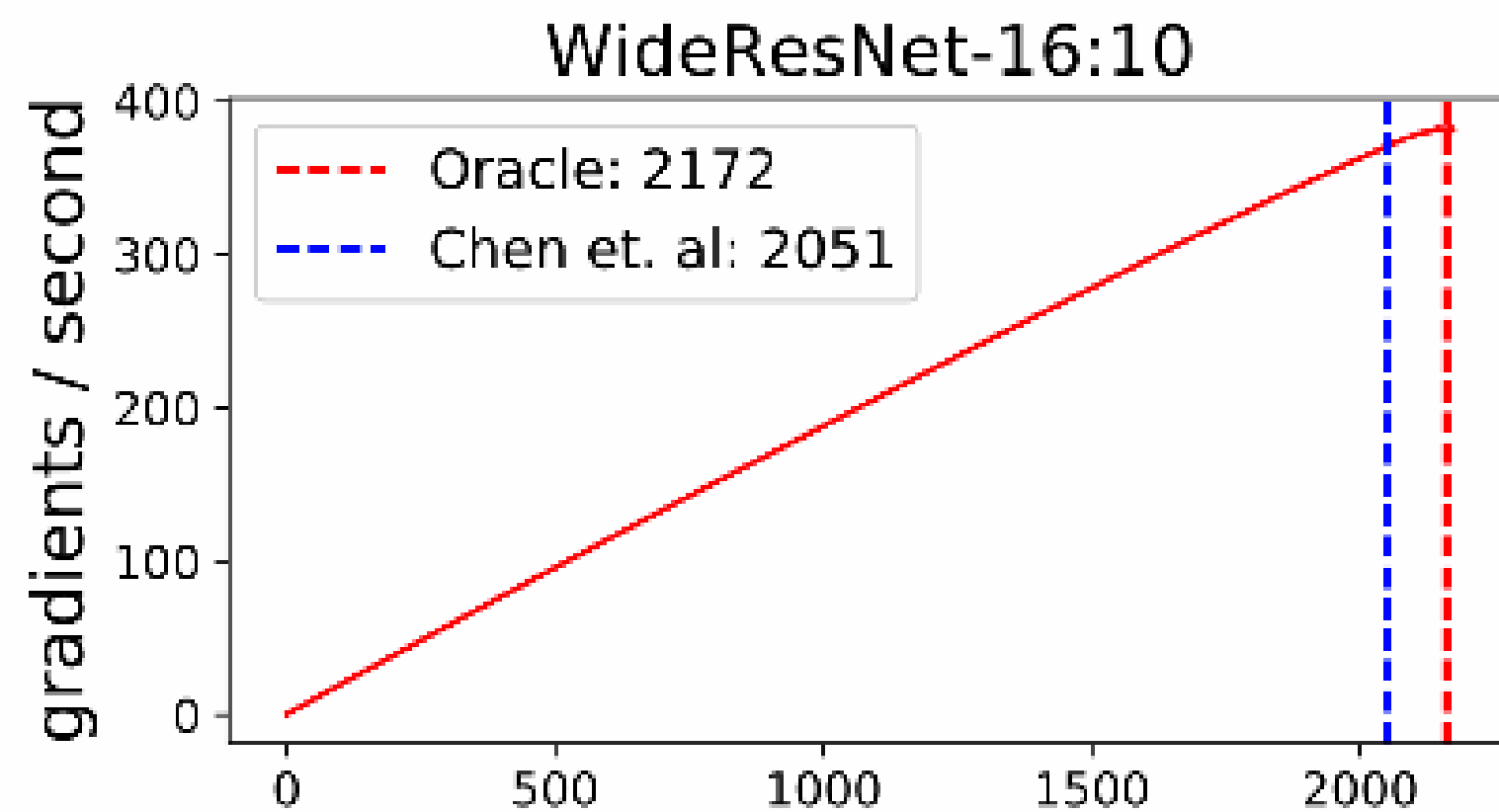# Async SGD or Stale Gradients

# Sync SGD or Straggled Workerks

# Why authers decided to write this paper

# Metric to optimize

$$\Omega(c) = \frac{c}{\tilde{x}_{(c)}}$$

$$\tilde{x}_{(1)}, \tilde{x}_{(2)}, \ldots, \tilde{x}_{(n)}$$

$$\arg\max_c \Omega(c) = \arg\max_c \mathbb{E}\left[\frac{c}{\tilde{x}_{(c)}}\right]$$

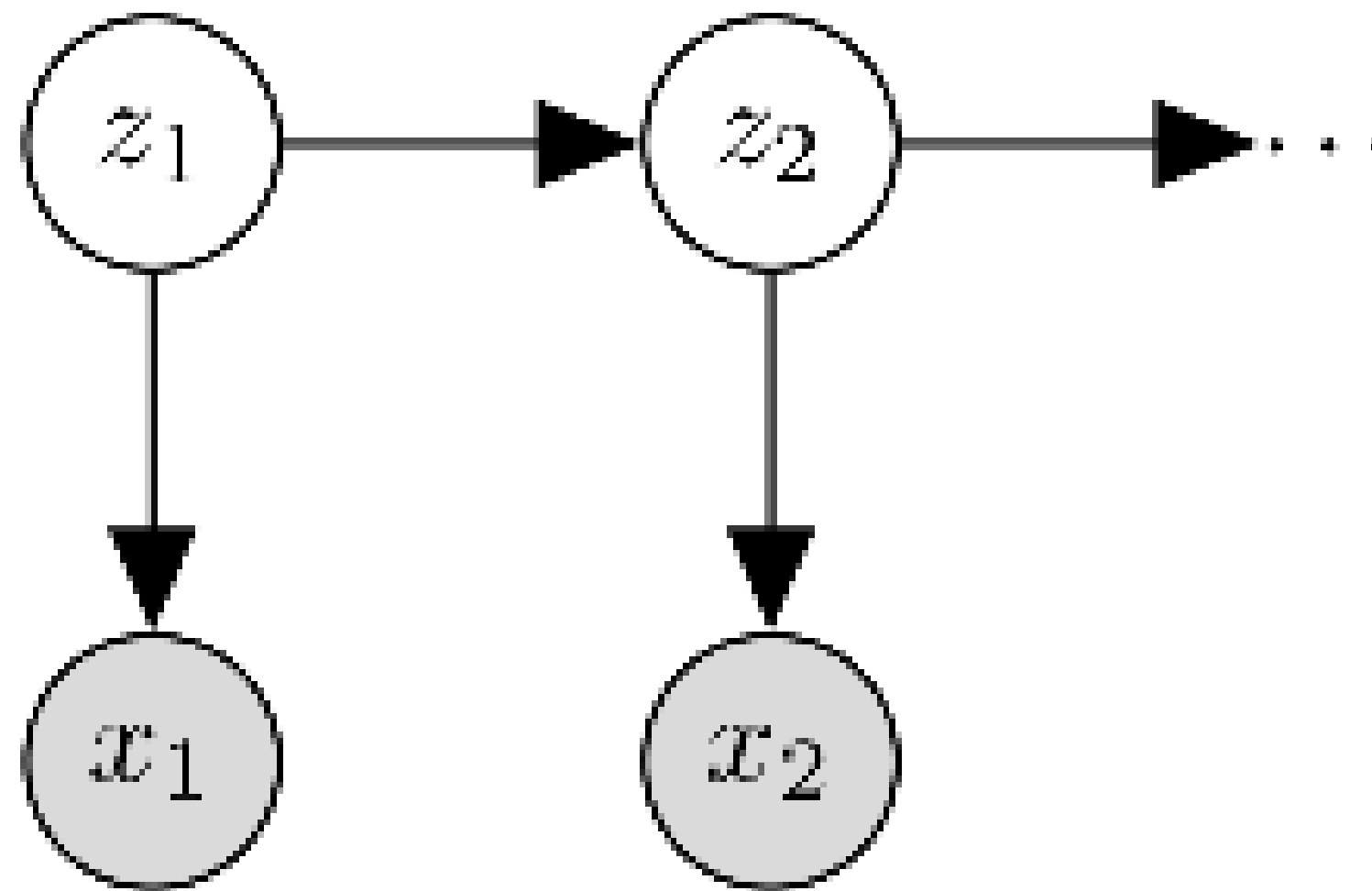# Elfving Cutoff

$$p(\tilde{x}_{(j)}) = Z(n,j) \int_{-\infty}^{\infty} x[\Phi(x)]^{j-1}[1 - \Phi(x)]^{n-j}p(x)dx$$

$$\mathbb{E}[\tilde{x}_{(j)}] \approx \mu_t + \Phi^{-1}\left(\frac{n - \frac{\pi}{8}}{j - \frac{\pi}{4} + 1}; 0, 1\right)\sigma_t$$
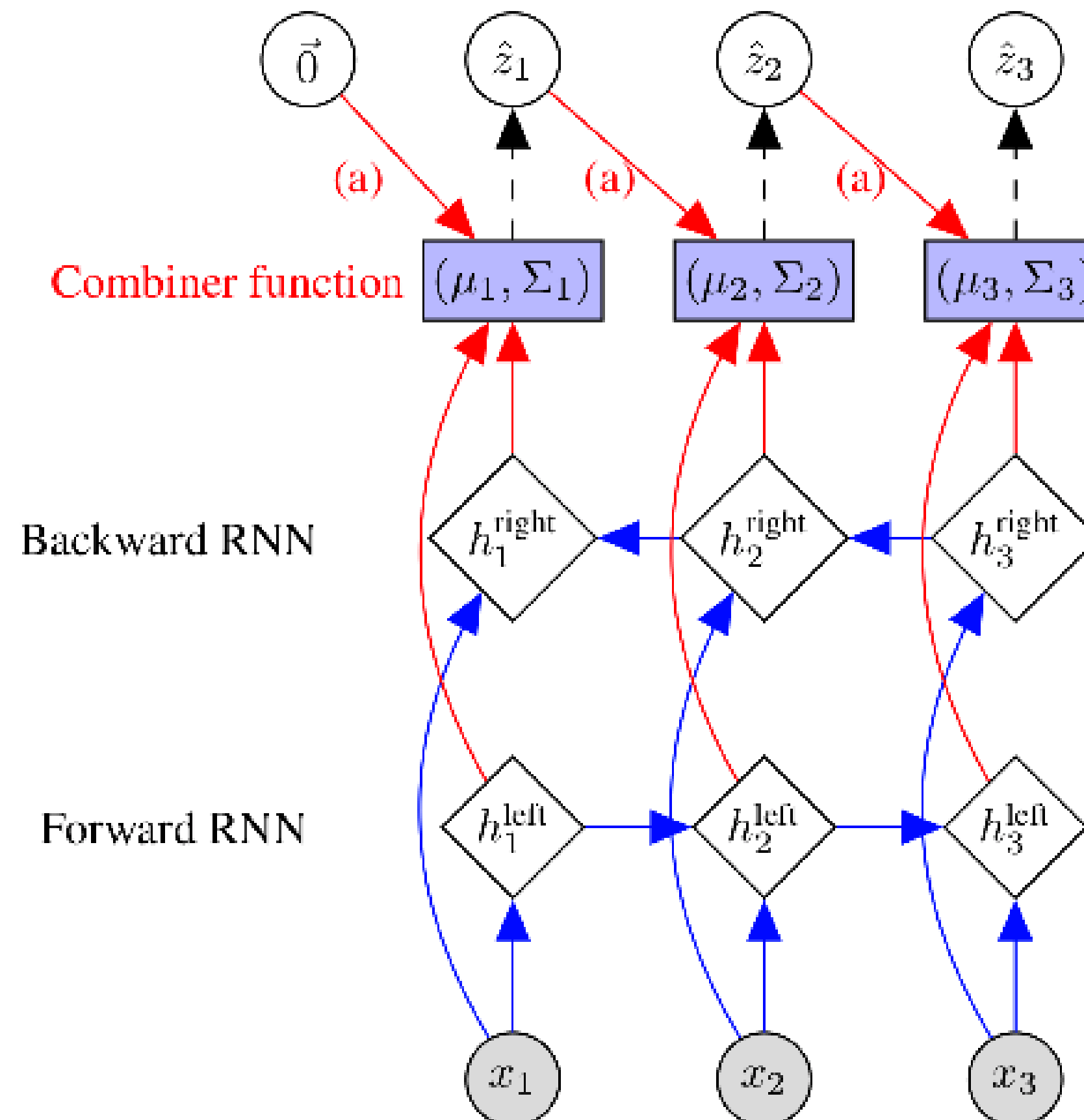
# Structured Inference Networks for Nonlinear Gaussian State Space

# Gaussian State Space and Deep Markov Models



$$z_t \sim \mathcal{N}(G_\alpha(z_{t-1}, \Delta_t), S_\beta(z_{t-1}, \Delta_t)) \quad \text{(Transition)}$$
$$x_t \sim \Pi(F_\kappa(z_t)) \quad \text{(Emission)}$$
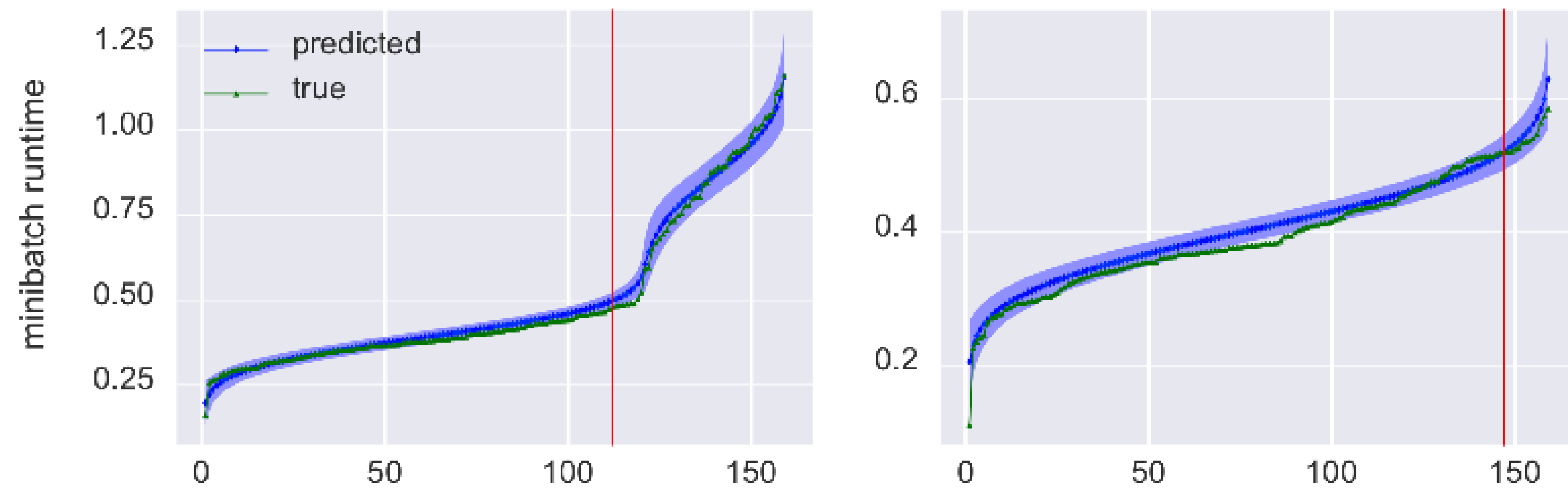
# Structured Inference Network

# Back to BDSGD

$$p_\theta(\boldsymbol{x}_{T-\ell:T}, \boldsymbol{z}_{T-\ell:T}) = \prod_{i=T-\ell}^{T} p_\theta(\boldsymbol{z}_i|\boldsymbol{z}_{i-1}) \prod_{i=T-\ell}^{T} p_\theta(\boldsymbol{x}_i|\boldsymbol{z}_i)$$

$$p(\boldsymbol{x}_{T+1}|\boldsymbol{x}_{T-\ell:T}) = \int p_\theta(\boldsymbol{x}_{T+1}|\boldsymbol{z}_{T+1})p_\theta(\boldsymbol{z}_{T+1}|\boldsymbol{z}_T)p(\boldsymbol{z}_{T-\ell:T}|\boldsymbol{x}_{T-\ell:T})d\boldsymbol{z}_{T-\ell:T+1}$$

$$\text{ELBO} = \mathbb{E}_{q_\phi(\boldsymbol{z}_{T-\ell:t}|\boldsymbol{x}_{T-\ell:T})} \log \left( \frac{p_\theta(\boldsymbol{x}_{T-\ell:t}, \boldsymbol{z}_{T-\ell:t})}{q_\phi(\boldsymbol{z}_{T-\ell:t}|\boldsymbol{x}_{T-\ell:T})} \right)$$

$$q_\phi(\boldsymbol{z}_{T-\ell:t}|\boldsymbol{x}_{T-\ell:T}) = \prod_{t=T-\ell}^{T} q_\phi(\boldsymbol{z}_t|\boldsymbol{z}_{T-\ell:t}, \boldsymbol{x}_{T-\ell:T}).$$
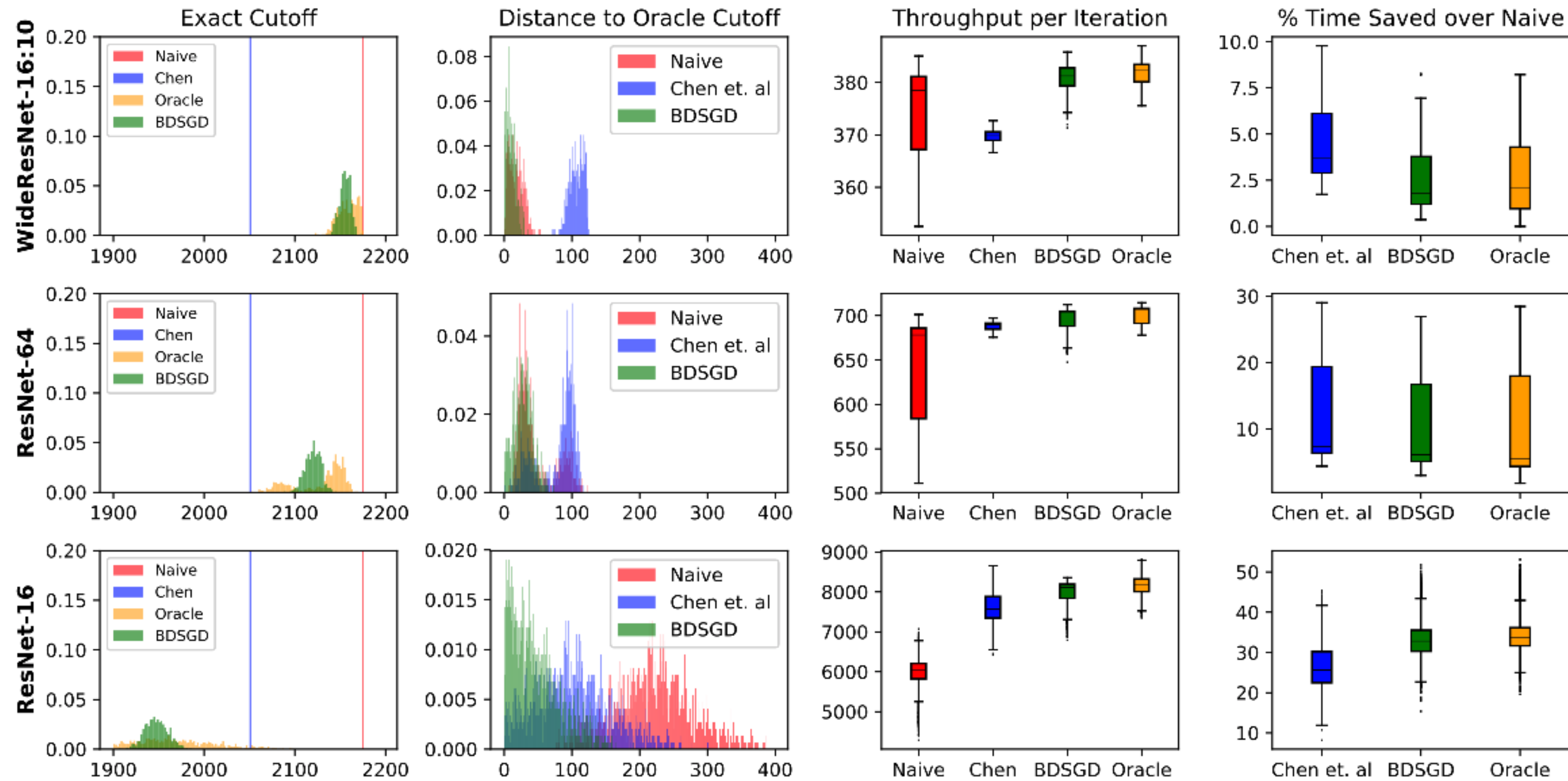
# Results



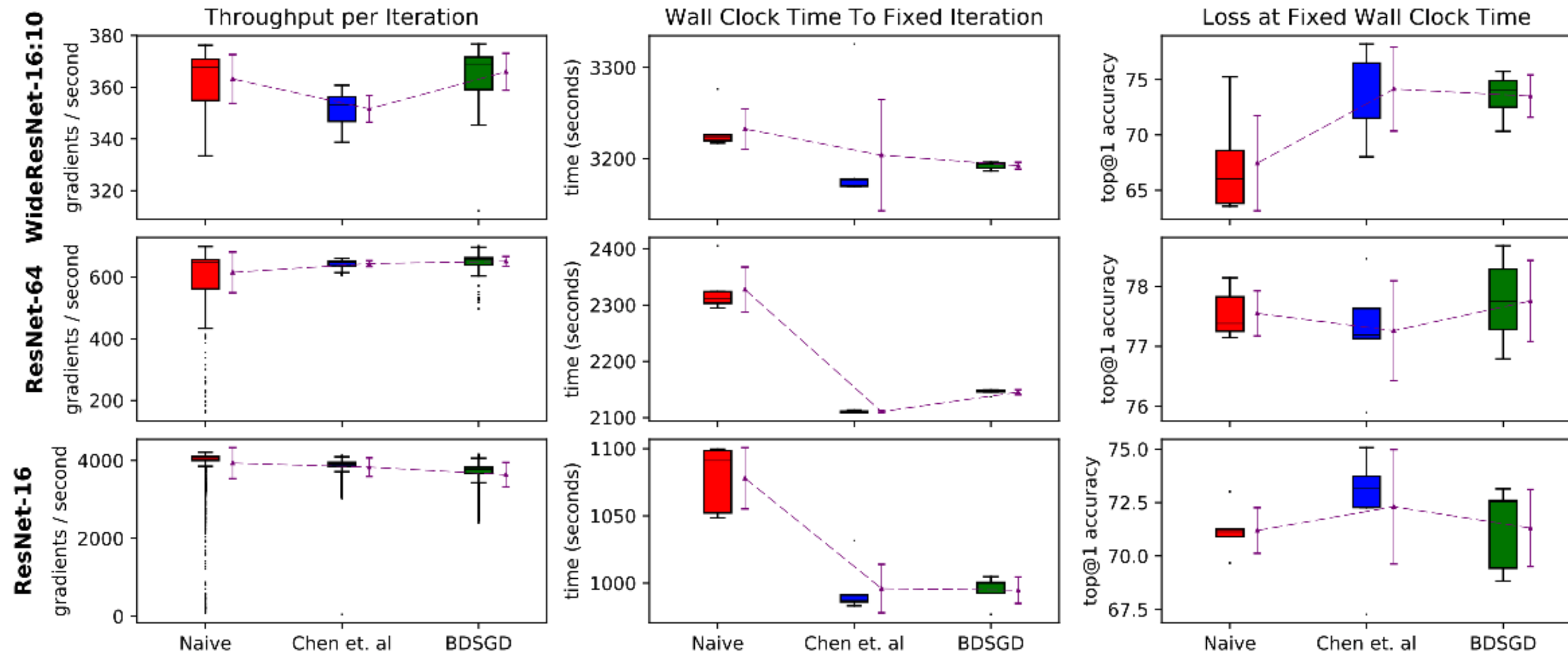(a) SGD runtime profiles for two iterations

# Results



(b) MNIST full training run

# Results

# Results

# Interesting Facts

VAE inference is also called *amortised inference*

Authers use DMM for predicting straggled
workers time

# Papers

**BDSGD**

https://papers.nips.cc/paper/7874-bayesian-distributed-stochastic-gradient-descent

**Structured Inference Network**

https://arxiv.org/abs/1609.09869