

# Deep Spectral Ranking

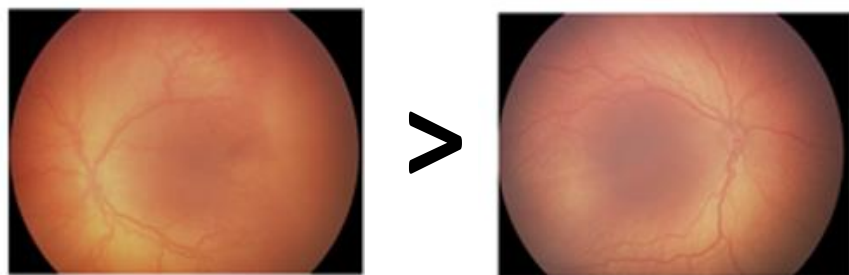
**Presenter: İlkey Yıldız, [yildizi@ece.neu.edu](mailto:yildizi@ece.neu.edu)**

İlkey Yıldız<sup>1</sup>, Jennifer Dy<sup>1</sup>, Deniz Erdoğan<sup>1</sup>, Susan Ostmo<sup>2</sup>, J. Peter Campbell<sup>2</sup>,  
Michael F. Chiang<sup>3</sup>, Stratis Ioannidis<sup>1</sup>

1. Electrical and Computer Engineering Dept., Northeastern University, Boston, MA, USA
2. Casey Eye Institute, Oregon Health and Science University, Portland, OR, USA
3. National Eye Institute, National Institutes of Health, Bethesda, MD, USA

# Motivation

- Ranking observations are classic in many domains.



*More diseased image*



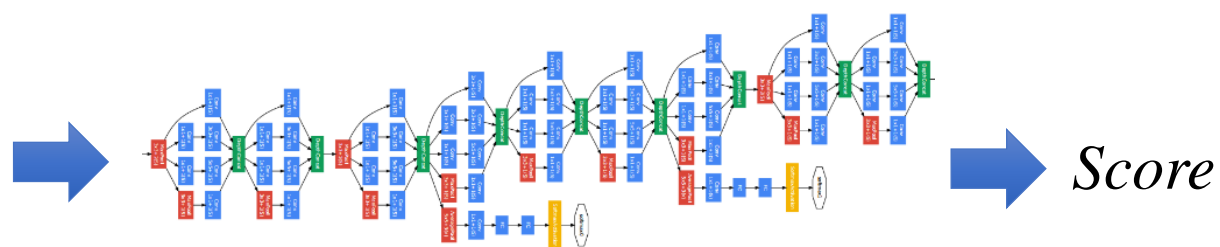
*Better Movie*

- Inference:** Plackett-Luce Model
  - Each sample has a positive *score*.
  - Probability that a sample is ranked higher is proportional to this score.

# Motivation

- **Regression**

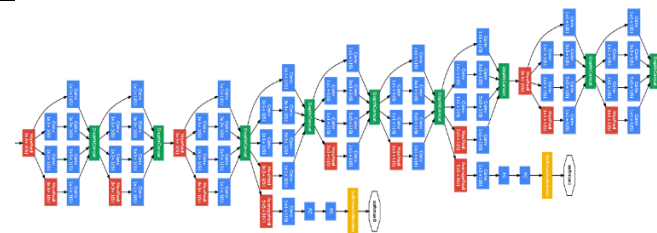
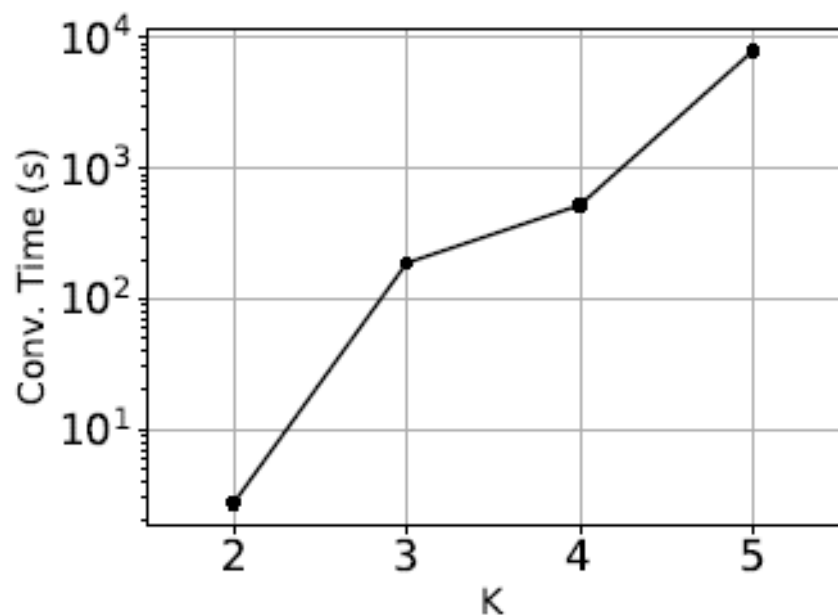
- **Shallow**: Newton's method [Tian et al., 2019]
- **DNN**: Siamese network attains 0.92 AUC from 80 samples! [Yıldız et al., 2019]



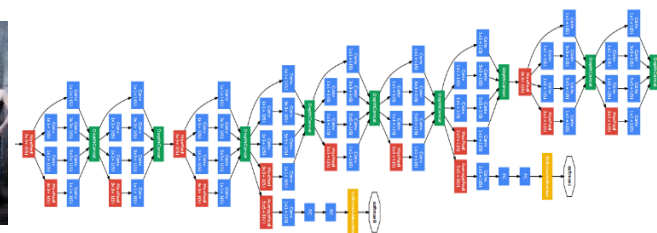
- [1] Tian, P., Guo, Y., ..., Chiang, M. F., Dy, J., Erdogmus, D., and Ioannidis, S. (2019). A severity score for Retinopathy of Prematurity. SIGKDD.
- [2] Yıldız, I., Tian, P., Dy, J., Erdogmus, D., Brown, J., ..., Chiang, M. F., and Ioannidis, S. (2019). Classification and comparison via neural networks. Neural Networks.

# Challenges

- Traditional DNN method: *siamese architecture* with  $K$  identical base networks
- Large –and potentially variable– memory footprint



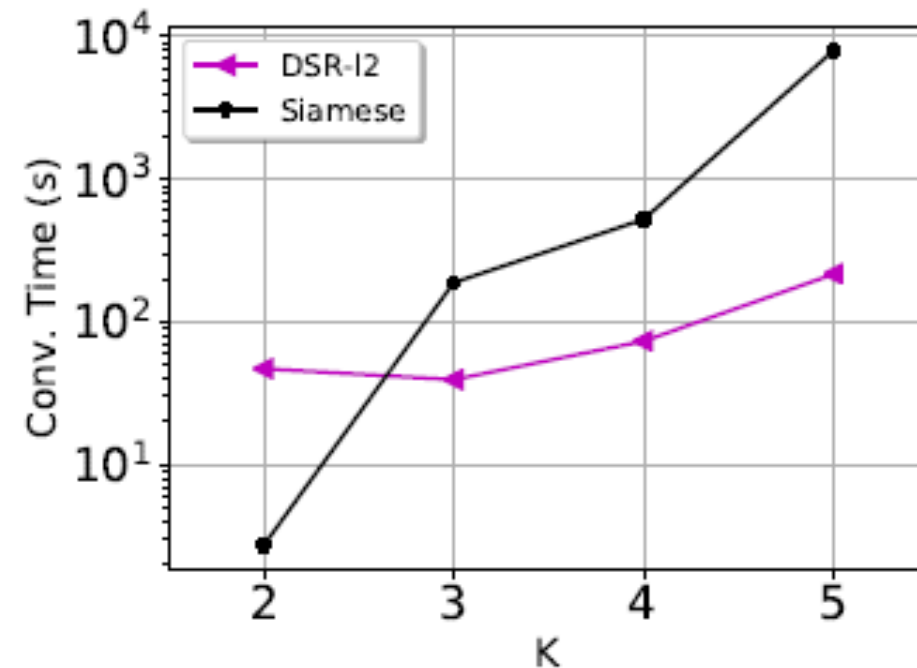
$K$  samples  $\Leftrightarrow$   
 $K$  networks



- An SGD epoch over  $\binom{n}{K} = O(n^K)$  observations
  - Exponential* training time

# Challenges

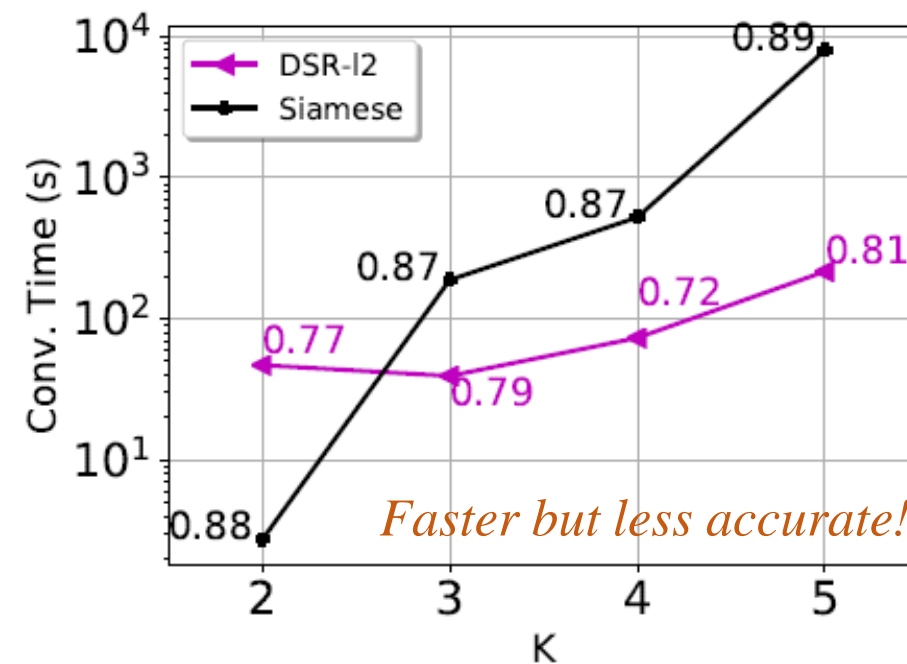
- A *spectral* algorithm for shallow regression
  - Solve via *ADMM* → Scores form the stationary distribution of a *Markov Chain*.
- Does not generalize well to deep models!
  - Scores form a *distribution*  $\leftrightarrow \ell_2$  penalty leads to vanishing gradients and stationary points!



[3] Yildiz, I., Dy, J., Erdogmus, D., ..., Chiang, M. F., and Ioannidis, S. (2020). Fast and accurate ranking regression. AISTATS.

# Challenges

- A *spectral* algorithm for shallow regression
  - Solve via *ADMM* → Scores form the stationary distribution of a *Markov Chain*.
- Does not generalize well to deep models!
  - Scores form a *distribution*  $\leftrightarrow \ell_2$  penalty leads to vanishing gradients and stationary points!



[3] Yildiz, I., Dy, J., Erdogmus, D., ..., Chiang, M. F., and Ioannidis, S. (2020). *Fast and accurate ranking regression*. AISTATS.

# Contributions

- Bridge the gap between *DNN* models and *spectral* algorithms for ranking regression
- Replace  $\ell_2$ -penalty of ADMM with *KL divergence*: still amenable to a spectral method!
- Significantly outperform  $\ell_2$ -penalty ADMM *and* siamese network

# Problem Formulation

## Plackett-Luce Model

- For each sample  $i \in \mathcal{N}$ , there exists a score  $\pi_i \in \mathbb{R}_+$ .
- Given query set of alternatives:  $A_\ell \subseteq \mathcal{N}$
- $m$  independent choice observations:  $c_\ell \in A_\ell, \ell \in \mathcal{M}$

$$\mathbf{P}(c_\ell | A_\ell, \boldsymbol{\pi}) = \pi_{c_\ell} / \sum_{j \in A_\ell} \pi_j = \pi_\ell / \sum_{j \in A_\ell} \pi_j$$

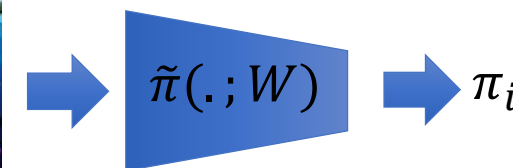


## Maximum Likelihood Estimation (MLE) for Ranking Regression

$$\begin{aligned} \text{Minimize}_{\boldsymbol{\pi}, \mathbf{W}} \quad & \mathcal{L}(\mathcal{D} | \boldsymbol{\pi}) \equiv \sum_{\ell=1}^m (\log \sum_{j \in A_\ell} \pi_j - \log \pi_\ell) \\ \text{subject to:} \quad & \boldsymbol{\pi} = \tilde{\boldsymbol{\pi}}(\mathbf{X}; \mathbf{W}) = [\tilde{\pi}_i = \tilde{\pi}(\mathbf{x}_i; \mathbf{W})]_{i \in \mathcal{N}} \\ & \boldsymbol{\pi} \geq \mathbf{0}, \end{aligned}$$



$\mathbf{x}_i$



- Traditional siamese architecture has base network:  $\tilde{\pi}(\cdot; W)$ .



# Alternating Direction Method of Multipliers (ADMM) with generalized penalty

$$\text{Minimize}_{\pi, W} \quad \mathcal{L}(\mathcal{D} | \pi) \equiv \sum_{\ell=1}^m (\log \sum_{j \in A_\ell} \pi_j - \log \pi_\ell)$$

subject to:  $\pi = \tilde{\pi}(X; W), \quad \pi \geq \mathbf{0},$

- **ADMM:** Decouple optimization of scores and parameters

$$L_\rho(\pi, W, \mathbf{y}) = \mathcal{L}(\mathcal{D} | \pi) + \mathbf{y}^\top (\pi - \tilde{\pi}(X; W)) + \rho \cdot D_p(\pi || \tilde{\pi}(X; W))$$

$$\pi^{k+1} = \arg \min_{\pi \in \mathbb{R}_+^n} L_\rho(\pi, W^k, \mathbf{y}^k)$$

Efficient spectral approach over the exponential ranking data!

$$W^{k+1} = \arg \min_{W \in \mathbb{R}^{d'}} \rho D_p(\pi^{k+1} || \tilde{\pi}(X; W)) - \mathbf{y}^k{}^\top \tilde{\pi}(X; W)$$

Linear in number of samples!

$$\mathbf{y}^{k+1} = \mathbf{y}^k + \rho(\pi^{k+1} - \tilde{\pi}(X; W^{k+1}))$$

Dual variable update

# Scores are still amenable to a spectral solution!

$$\frac{\partial L_\rho(\pi, W^k, y^k)}{\partial \pi_i} = 0 \quad \forall i \iff \sum_{j \neq i} \pi_j \lambda_{ji}(\pi) - \sum_{j \neq i} \pi_i \lambda_{ij}(\pi) = \pi_i \sigma_i, \quad (15)$$

**Theorem 4.2.** Eq. (15) are the balance equations of a continuous-time MC with transition rates:

$$\mu_{ji}(\pi) = \begin{cases} \lambda_{ji}(\pi) + \frac{2\pi_i \sigma_i \sigma_j}{\sum_{t \in \mathcal{N}_-} \pi_t \sigma_t - \sum_{t \in \mathcal{N}_+} \pi_t \sigma_t} & \text{if } j \in \mathcal{N}_+ \text{ and } i \in \mathcal{N}_- \\ \lambda_{ji}(\pi) & \text{otherwise,} \end{cases} \quad (16)$$

Stationary scores are also the stationary distribution of the continuous time MC.

$$\sigma_i(\pi) = \rho \frac{\partial D_p(\pi || \tilde{\pi}^k)}{\partial \pi_i} + y_i^k$$

**ILSRX:**  $\pi^{l+1} = \text{ssd}(M(\pi^l))$

# ADMM + KL proximal penalty

$l_2$  proximal penalty:  $D_p(\pi || \tilde{\pi}) = \|\tilde{\pi} - \pi\|_2^2$

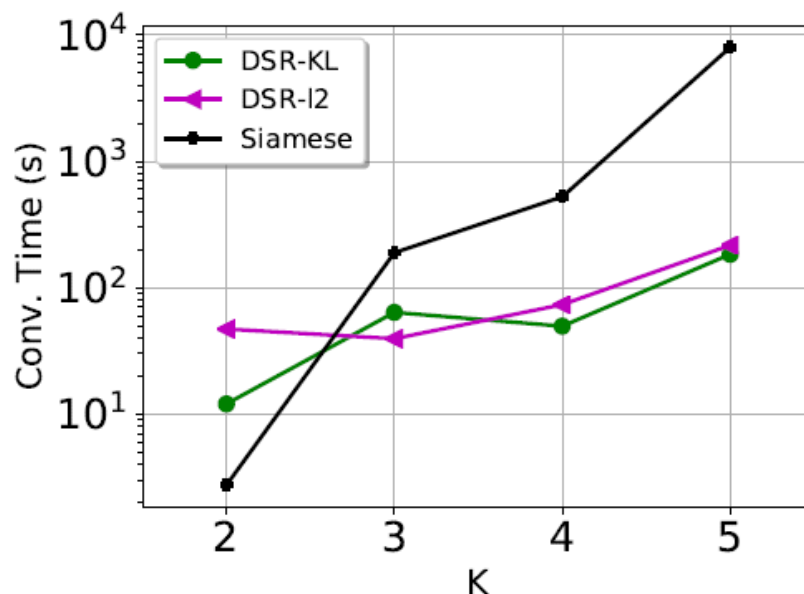
$$W^{k+1} = \arg \min_{W \in \mathbb{R}^{d'}} \|\tilde{\pi}(X; W) - (\pi^{k+1} + \frac{1}{\rho} y^k)\|_2^2$$

KL proximal penalty:  $D_p(\pi || \tilde{\pi}) = \sum_{i=1}^n \pi_i \log \frac{\pi_i}{\tilde{\pi}_i}$

$$W^{k+1} = \arg \min_{W \in \mathbb{R}^{d'}} \sum_{i=1}^n (-\frac{y_i^k}{\rho} \tilde{\pi}_i - \pi_i^{k+1} \log \tilde{\pi}_i)$$

- $l_2$  loss: Prone to vanishing gradients and reaching stationary points.

- Naturally suited as  $\pi = [\pi_i]_{i \in [n]} \in \mathbb{R}_+^n$  is a distribution.
- Max-entropy loss: better fitting and convergence



# ADMM + KL penalty

$l_2$  proximal penalty:  $D_p(\pi || \tilde{\pi}) = \|\tilde{\pi} - \pi\|_2^2$

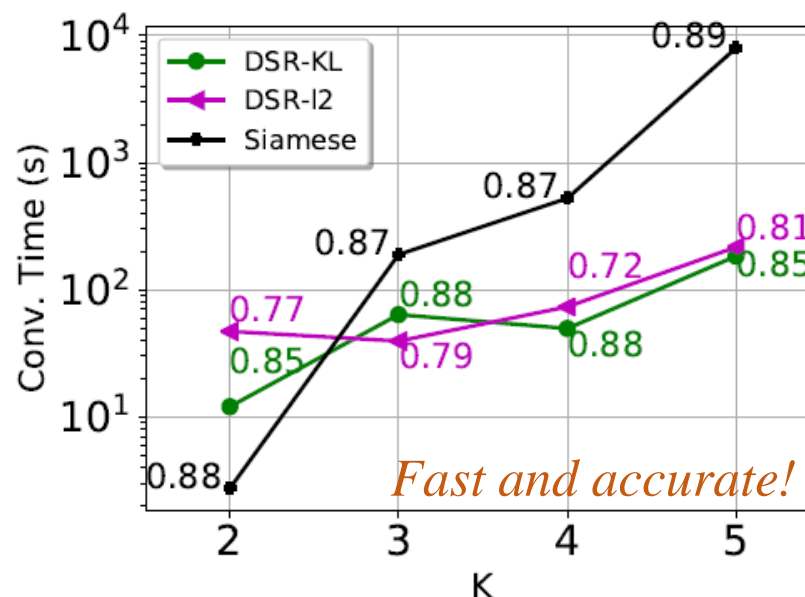
$$W^{k+1} = \arg \min_{W \in \mathbb{R}^{d'}} \|\tilde{\pi}(X; W) - (\pi^{k+1} + \frac{1}{\rho} y^k)\|_2^2$$

KL proximal penalty:  $D_p(\pi || \tilde{\pi}) = \sum_{i=1}^n \pi_i \log \frac{\pi_i}{\tilde{\pi}_i}$

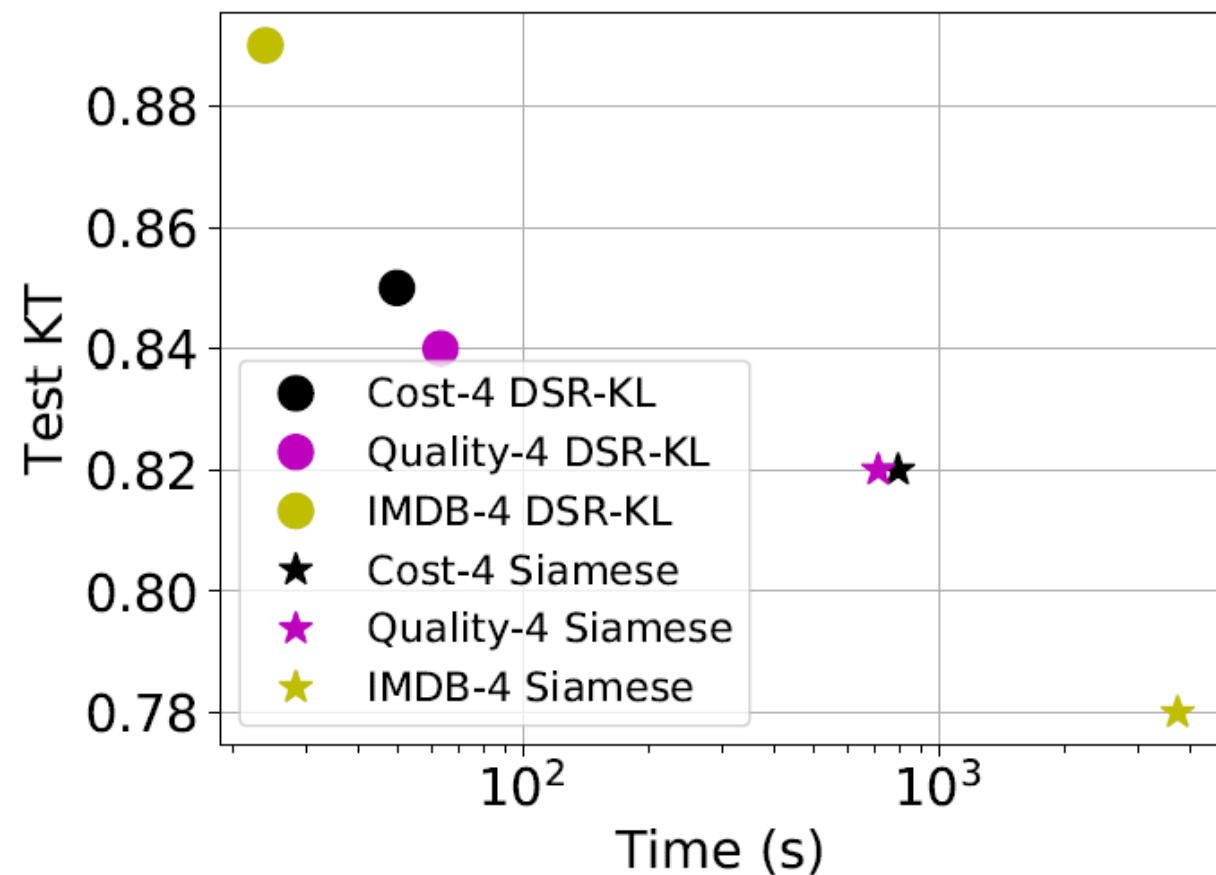
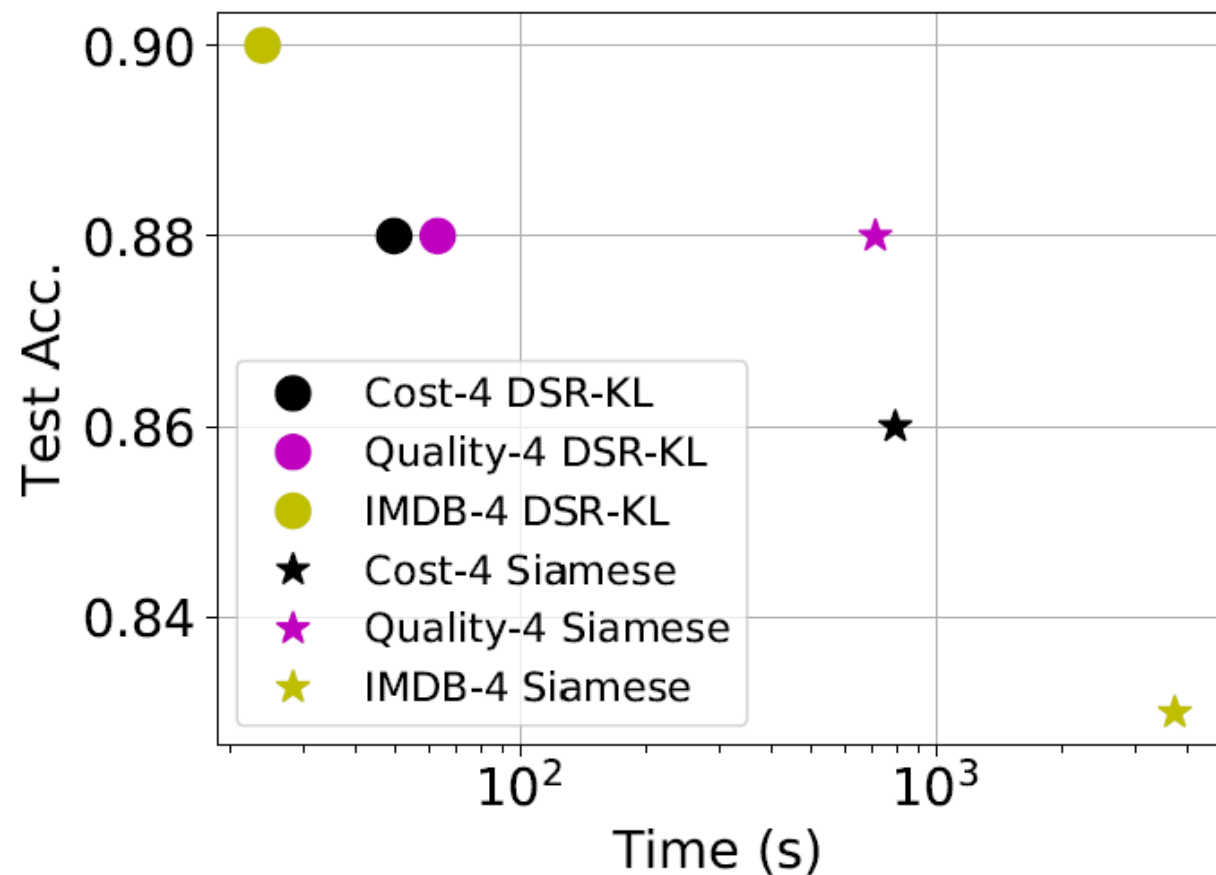
$$W^{k+1} = \arg \min_{W \in \mathbb{R}^{d'}} \sum_{i=1}^n (-\frac{y_i^k}{\rho} \tilde{\pi}_i - \pi_i^{k+1} \log \tilde{\pi}_i)$$

- $l_2$  loss: Prone to vanishing gradients and reaching stationary points.

- Naturally suited as  $\pi = [\pi_i]_{i \in [n]} \in \mathbb{R}_+^n$  is a distribution.
- Max-entropy loss: better fitting and convergence



# Faster Training AND Better Predictions!



# THANK YOU!

Presenter: İlkay Yıldız, [yildizi@ece.neu.edu](mailto:yildizi@ece.neu.edu)



Supported by NIH (R01EY019474), NSF (SCH-1622542 at MGH, SCH-1622536 at Northeastern, SCH-1622679 at OHSU), Facebook Statistics Research Award, and by unrestricted departmental funding from Research to Prevent Blindness (OHSU)