

Vector compression for similarity search using Multi-layer Sparse Ternary Codes

Sohrab Ferdowsi, Slava Voloshynovskiy, Dimche Kostadinov

Department of Computer Science, University of Geneva, Switzerland

{Sohrab.Ferdowsi, svolos, Dimche.Kostadinov}@unige.ch

http://sip.unige.ch

Overview

Applications:

- Large-scale retrieval systems
- Learned compression of feature vectors
- Compressed representation useful for fast similarity search

Contributions:

- Rate-distortion (R-D) study of ternary and binary encoding
- Designing R-D efficient multi-layer Sparse Ternary Codes (STC)

Background

- Ability to **search for similarity** within a database is crucial for modern retrieval systems.
- A wide-spread solution is binary hashing.
- We proposed ternary hashing [WIFS'16] as an alternative to binary hashing.
- We showed that ternary encoding has higher **coding gain** than binary encoding [ISIT'17].
- Here we extend ternary encoding for the task of compression, so that we can have list-refinement.
- Our design challenge: To have good R-D performance within STC limitations.

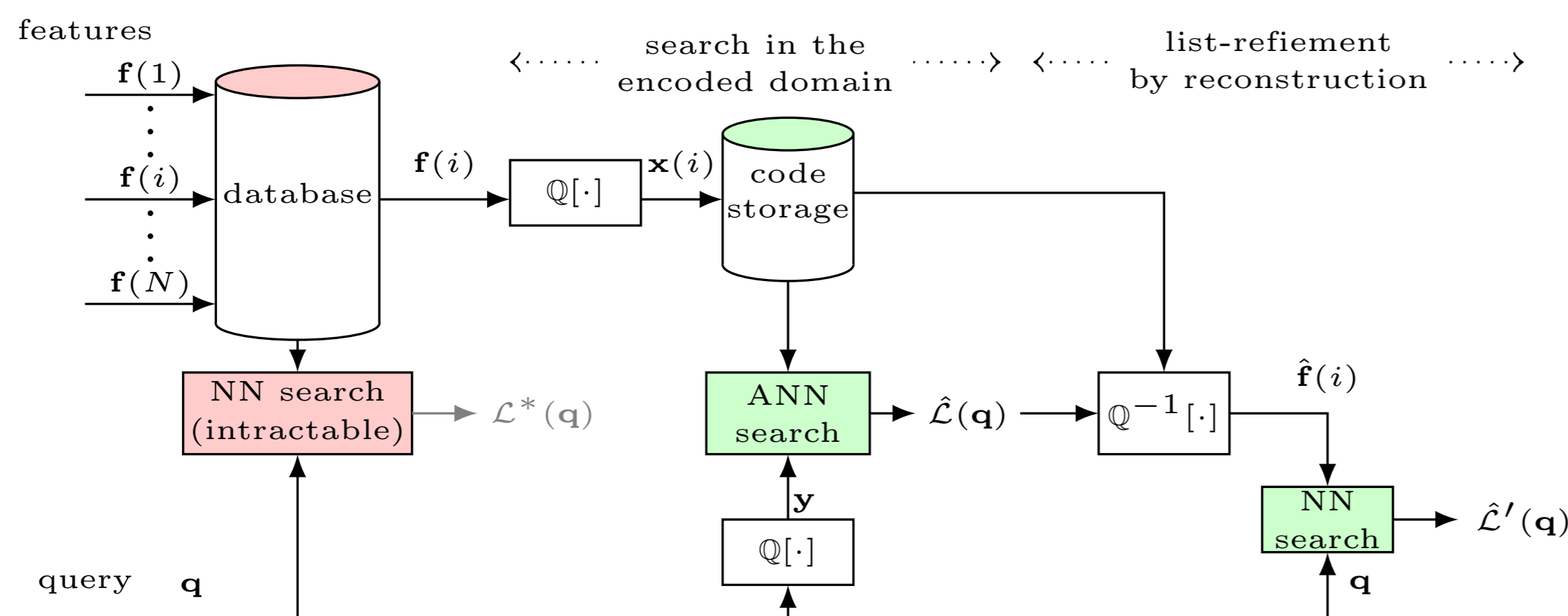
Problem formulation: ANN search

Similarity search:

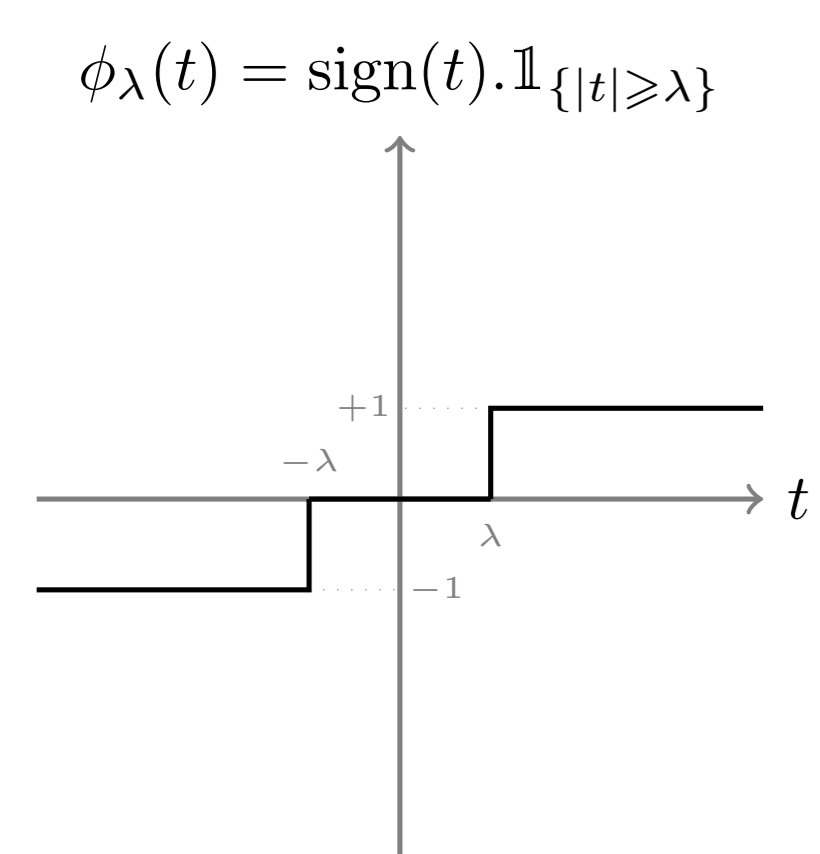
- (Exact) Nearest Neighbor (NN) search: $\mathcal{L}(\mathbf{q}) = \{1 \leq i \leq N | d_E(\mathbf{f}(i), \mathbf{q}) \leq \epsilon\}$
- Approximate Nearest Neighbor (ANN) search: $\hat{\mathcal{L}}(\mathbf{q}) = \{1 \leq i \leq N | d_H(\mathbb{Q}[\mathbf{f}(i)], \mathbb{Q}[\mathbf{q}]) \leq \epsilon\}$
- Our solution: List-refinement with reconstruction $\hat{\mathcal{L}}'(\mathbf{q}) = \{i \in \hat{\mathcal{L}}(\mathbf{q}) | d_E(\mathbb{Q}^{-1}[\mathbb{Q}[\mathbf{f}(i)]], \mathbf{q}) \leq \epsilon\}$

Compression:

- Encoding: $\mathbf{x} = \mathbb{Q}[\mathbf{f}]$
- Reconstruction: $\hat{\mathbf{f}} = \mathbb{Q}^{-1}[\mathbf{x}]$
- Rate: $\mathcal{R} = \frac{1}{n} \mathbb{E}[\# \text{ bits used to represent } \mathbf{x}]$
- Distortion: $\mathcal{D} = \mathbb{E}[d_E(\mathbf{F}, \hat{\mathbf{F}})]$
- $d_E(\mathbf{a}, \mathbf{b}) \triangleq \frac{1}{n} \|\mathbf{a} - \mathbf{b}\|_2^2$



Single-layer Sparse Ternary Codes (STC)



Encoding:

$$\mathbf{x} = \phi_\lambda(\mathbf{A}\mathbf{f}) \odot \beta$$

(projection + ternarization + re-weighting)

Reconstruction:

$$\hat{\mathbf{f}} = \mathbf{B}\mathbf{x} = \mathbf{B}\phi_\lambda(\mathbf{A}\mathbf{f}) \odot \beta$$

(projection)

Optimizing single-layer STC

Back-projection B:

- Decompose $\mathbf{B} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{B}'$, optimize \mathbf{B}' :

$$\mathbf{B}' = \underset{\mathbf{B}'}{\operatorname{argmin}} \|\mathbf{F} - (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{B}' \mathbf{X}\|_2^2$$

$$= \underset{\mathbf{B}'}{\operatorname{argmin}} \operatorname{Tr} \left[-2\mathbf{A}\mathbf{A}^T \mathbf{A}\mathbf{F}\mathbf{X}^T \mathbf{B}'^T + \mathbf{B}' \mathbf{X}\mathbf{X}^T \mathbf{B}'^T \mathbf{A}\mathbf{A}^T \right]$$

$$\Rightarrow \mathbf{B}' = \mathbf{A}\mathbf{F}\mathbf{X}^T (\mathbf{X}\mathbf{X}^T)^{-1}$$

Projection A:

- To have un-correlated \mathbf{X} :

$$\mathbf{C}_F \triangleq \frac{1}{n} \mathbb{E}[\mathbf{F}\mathbf{F}^T] = \mathbf{U}_F \Sigma_F \mathbf{U}_F^T$$

- Simply as in PCA: $\mathbf{A} = \mathbf{U}_F^T$
- $\tilde{\mathbf{X}} \triangleq \mathbf{A}\mathbf{F} \sim \mathcal{N}(\mathbf{0}, \Sigma_F)$
- With this choice of \mathbf{A} : $\Rightarrow \mathbf{B}' = \mathbf{I}_n, \mathbf{B} = \mathbf{A}^T$

$$\mathcal{D} = \mathbb{E}[d_E(\mathbf{F}, \hat{\mathbf{F}})] = \frac{1}{n} \mathbb{E}[\|\mathbf{F} - \mathbf{A}^T \mathbf{X}\|_2^2] = \frac{1}{n} \mathbb{E}[\|\mathbf{A}\mathbf{F} - \mathbf{X}\|_2^2] = \frac{1}{n} \mathbb{E}[\|\tilde{\mathbf{X}} - \phi_\lambda(\tilde{\mathbf{X}}) \odot \beta\|_2^2]$$

- Distortion per each dimension ($\mathcal{D} = \sum_{i=1}^n \mathcal{D}_i$):

$$\mathcal{D}_i = \mathbb{E}[(\tilde{x}_i - \beta_i \phi_\lambda(\tilde{x}_i))^2] = \int_{-\infty}^{-\lambda} (\tilde{x}_i + \beta_i)^2 p(\tilde{x}_i) d\tilde{x}_i + \int_{-\lambda}^{+\lambda} \tilde{x}_i^2 p(\tilde{x}_i) d\tilde{x}_i + \int_{+\lambda}^{+\infty} (\tilde{x}_i - \beta_i)^2 p(\tilde{x}_i) d\tilde{x}_i$$

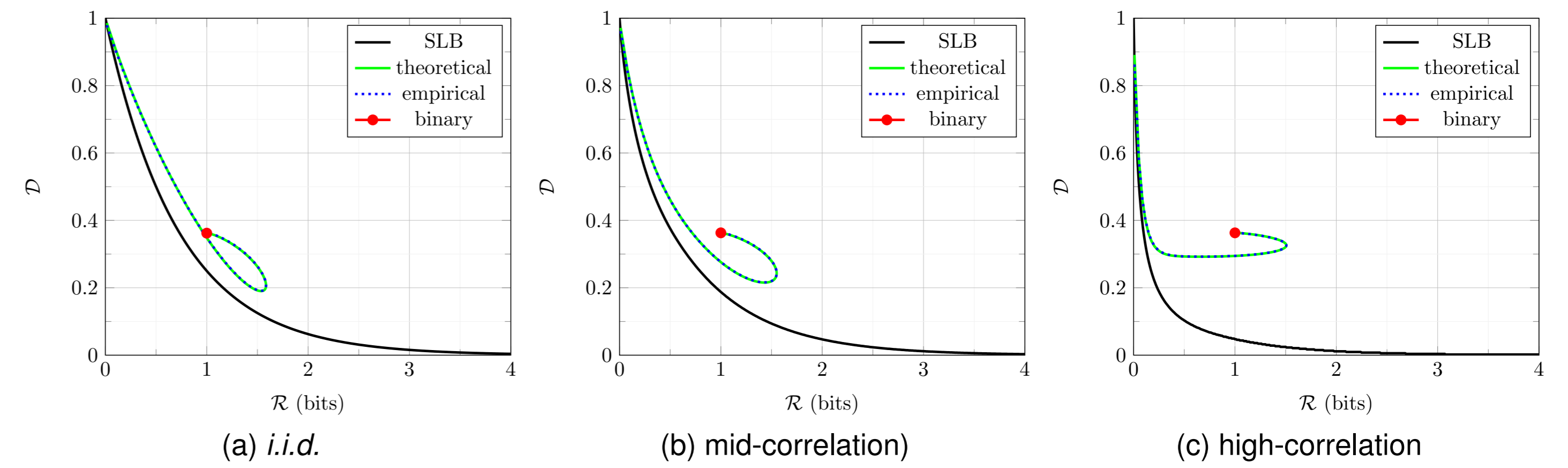
$$\Rightarrow \mathcal{D}_i = \sigma_i^2 + 2\beta_i^2 \mathcal{Q}\left(\frac{\lambda}{\sigma_i}\right) - \frac{4\beta_i \sigma_i}{\sqrt{2\pi}} \exp\left(-\frac{\lambda^2}{2\sigma_i^2}\right)$$

- Optimal Re-weighting vector: $\beta_i^* = \underset{\beta_i}{\operatorname{argmin}} \mathcal{D}_i = \frac{\sigma_i \exp\left(-\frac{\lambda^2}{2\sigma_i^2}\right)}{\sqrt{2\pi} \mathcal{Q}\left(\frac{\lambda}{\sigma_i}\right)}$

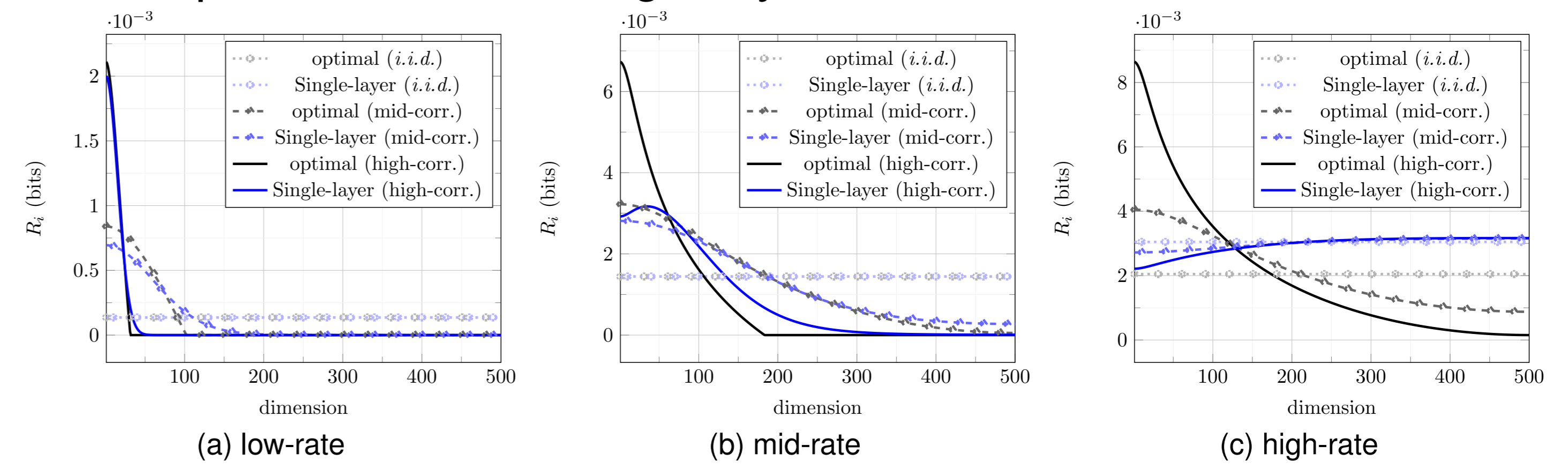
- Rate: $\mathcal{R} = \frac{1}{n} H(\mathbf{X}) = \frac{1}{n} \sum_{i=1}^n H(X_i) = -\frac{1}{n} \sum_{i=1}^n (2\alpha_i \log_2(\alpha_i) + (1 - 2\alpha_i) \log_2(1 - 2\alpha_i))$

- Sparsity per each dimension: $\alpha_i \triangleq \mathbb{P}[X_i = +\beta_i] = \mathbb{P}[X_i = -\beta_i]$

R-D performance of single-layer STC on AR(1) Gaussian sources

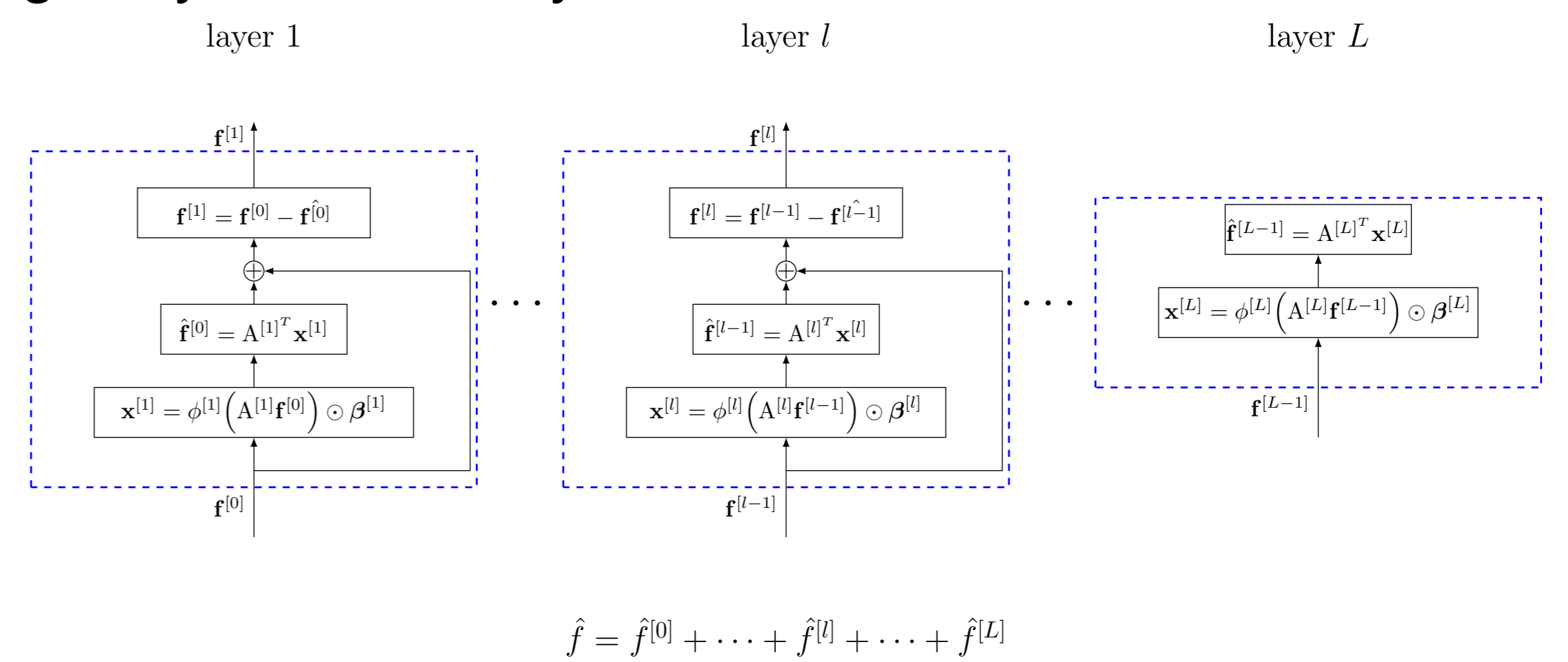


Poor R-D performance for single-layer: rate mismatch

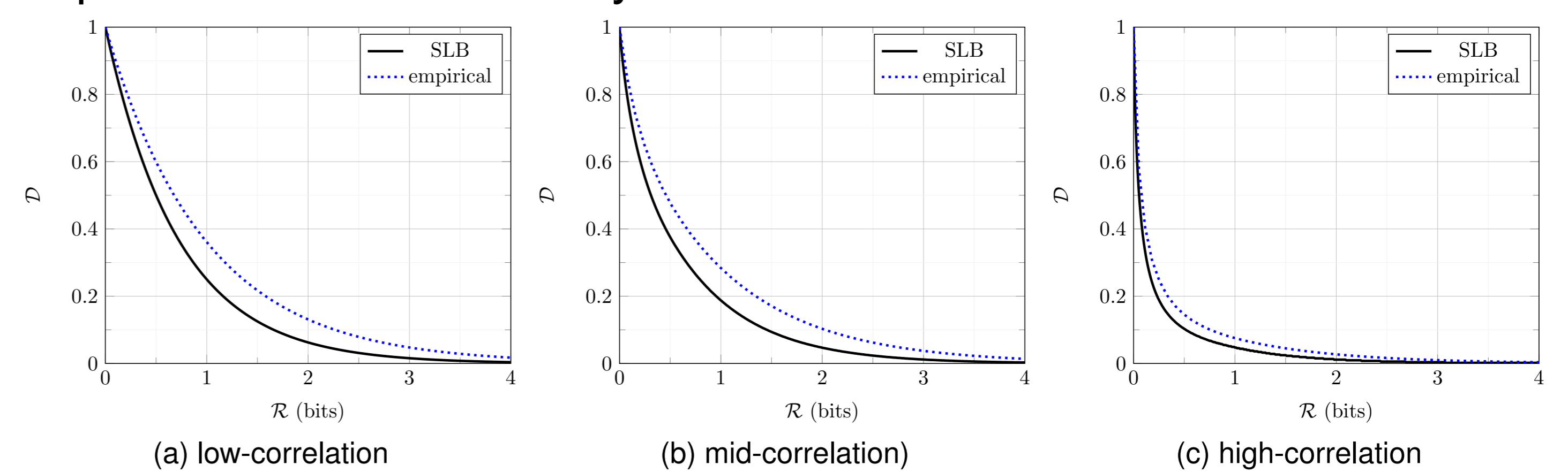


- Optimal rate allocation is calculated using the "reverse water-filling" paradigm from information theory.
- At higher rates, rate allocation deviates largely from optimal assignment.
- Binary encoding is a special case of ternary encoding with zero sparsity and hence rate is very high.

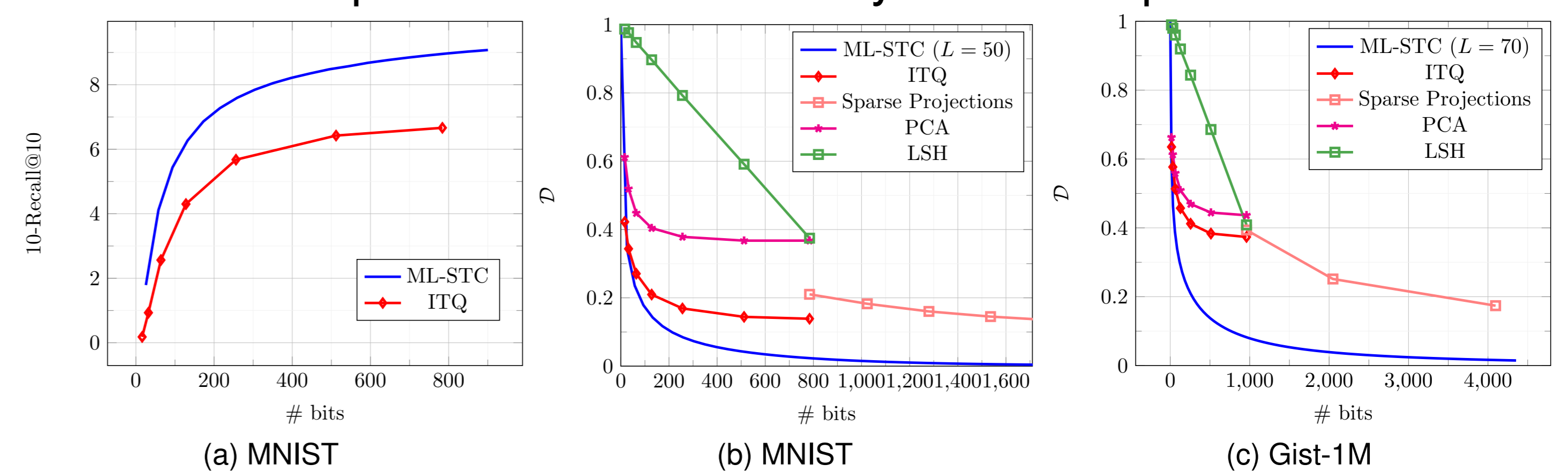
From single-layer to multi-layer architecture



R-D performance for multi-layer STC



Search and R-D performance of multi-layer STC on public databases



Conclusions

- Single-layer encoding is insufficient to provide good R-D performance at high rates.
- Residual-based multi-layer encoding can provide reasonable R-D performance.
- Since binary-encoding has rate mismatch, it cannot benefit from multi-layer encoding.
- Ternary encoding with high sparsity has low rate mismatch and can benefit from multi-layer encoding.
- Future work: Joint learning of all layers.
- Python implementation: <https://github.com/sssohrab/DSW2018>