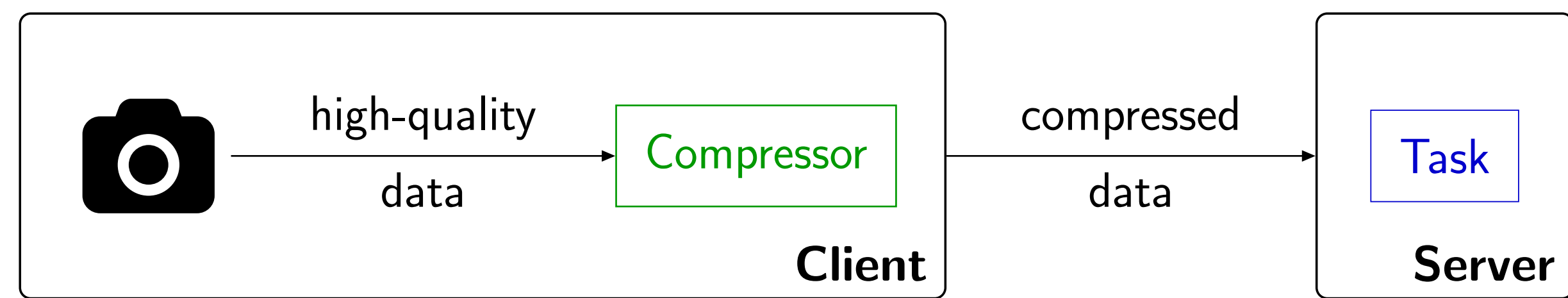


Introduction



Motivation: a resource constrained *client* offloads costly task-related computations to a remote *server* (edge/cloud computing).

Open need: design task-aware source coding schemes which provides *effective* representations of the source data.

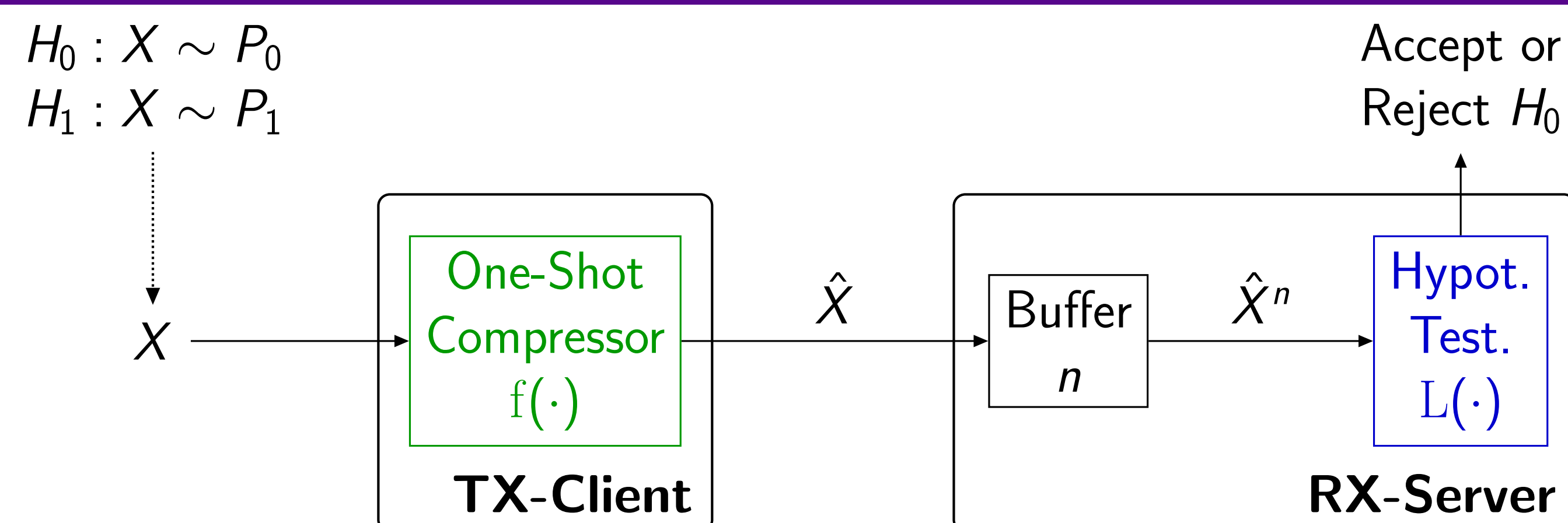
Assumptions:

- ▶ task: binary hypothesis testing;
- ▶ client: constrained device which cannot perform task locally, does not have memory and can only do simple scalar compression;
- ▶ server: hypothesis testing on a block of compressed samples.

Our work: single-shot fixed-length compression for hypothesis testing.

- ▶ problem formulation;
- ▶ analyze the error performance;
- ▶ propose a task-oriented compression algorithm for hypothesis testing.

System Model



Source	Compressor	Hypothesis Testing
$x \in \mathcal{X} = \{1, \dots, \mathcal{X} \}$	$f: \mathcal{X} \rightarrow \mathcal{M} = \{1, \dots, M\}$	$L(\hat{X}^n) \stackrel{\hat{\theta}=0}{\geq} \log T \stackrel{\hat{\theta}=1}{\leq}$
$X \sim P_\theta(x), \theta \in \{0, 1\}$	$\hat{X} = f(X), \hat{X} \sim \hat{P}_\theta(\hat{X})$	

Fixed rate compression $R = \log M$. We consider $M < |\mathcal{X}|$.

Task: binary hypothesis testing.

- ▶ if type-I error $< \epsilon$, then type-II error β_n^ϵ (accept H_0 when H_1 is true) decays exponentially in n as $\gamma = -\lim_{n \rightarrow \infty} \frac{1}{n} \log \beta_n^\epsilon$;
- ▶ **our performance metric:** type-II error exponent γ ;
- ▶ Chernoff-Stein [1]: optimal type-II error exponent is $\gamma^* = D(P_0||P_1)$ when there is no compression;
- ▶ with compression: error exponent depends on (f, R) : $\gamma_f(R)$;
- ▶ **compression penalty:** $\Delta_f(R) = D(P_0||P_1) - \gamma_f(R)$.

Hypothesis Testing under Single-shot Compression

Hypothesis test on $\hat{X} \sim \hat{P}_\theta$:

- ▶ log-likelihood ratio test on \hat{X}^n is optimal;
 - ▶ optimal error exponent is $\gamma_f^*(R) = D(\hat{P}_0||\hat{P}_1)$.
- \implies **Compression penalty:** $\Delta_f(R) = D(P_0||P_1) - D(\hat{P}_0||\hat{P}_1)$

Proposition 1. Expression for $\Delta_f \geq 0$:

$$\Delta_f = \sum_{\hat{x}=1}^M \hat{P}_0(\hat{x}) D(P_0(x|\hat{x}) || P_1(x|\hat{x}))$$

$P_\theta(X|\hat{X}) = \frac{P_\theta(X)}{\hat{P}_\theta(\hat{X})} \mathbb{1}\{\hat{X} = f(X)\}$ is the posterior of X given $\hat{X} = f(X)$.

- ▶ Note that a good task-aware compression strategy combines X that have similar posteriors $P_\theta(X|\hat{X})$.

Optimal compressor:

- ▶ $f^* = \arg \max_f D(\hat{P}_0||\hat{P}_1) = \arg \min_f \Delta_f$ s.t. $|f| \leq M$;
- ▶ optimization over each possible f , which induces a partition of M sets over \mathcal{X} (NP-hard).

Proposed Compressor Scheme

Optimal one-step compression from $|\mathcal{X}|$ to $|\mathcal{X}| - 1$:

- ▶ f combines $\{a, b\} \subset \mathcal{X}$ and the others $x \in \mathcal{X} \setminus \{a, b\}$ are one-to-one;
- ▶ i.e., $f(a) = f(b) = m \in \mathcal{M}$, $f(i) = i \in \mathcal{M} \setminus \{m\}$;

Then,

$$f^* = \arg \min_{\{a,b\} \subset \mathcal{X}: f(a)=f(b)=m} \left\{ \hat{P}_0(m) D(P_0(x|m) || P_1(x|m)) \right\}. \quad (1)$$

Our “KL-greedy” compressor:

- ▶ iteratively reduce the alphabet size by 1 at each step, until the compressed alphabet has size M ;
- ▶ at each step, combine $\{a, b\}$ which minimize (1);
- ▶ note that this compressor can be determined in polynomial time.

Results

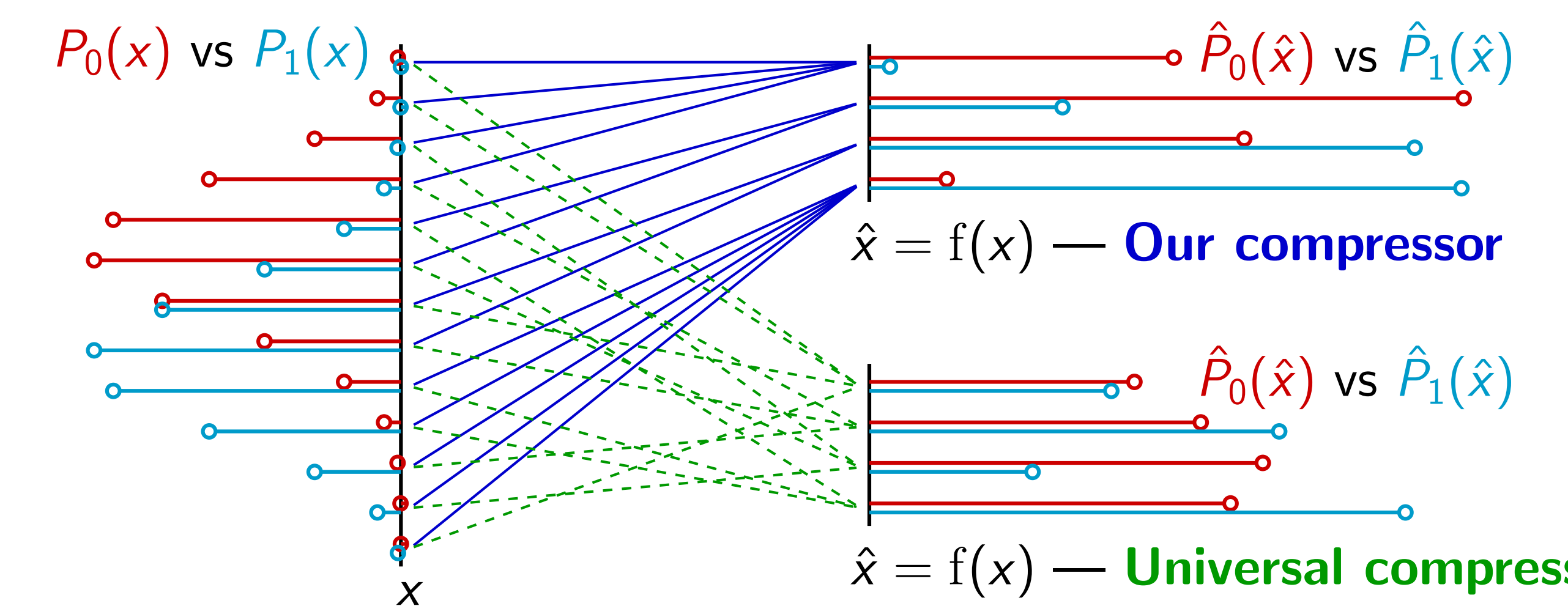
P_θ are shifted binomial distributions with different parameters.

Compare compression penalty Δ_f and empirical type-II error rate for:

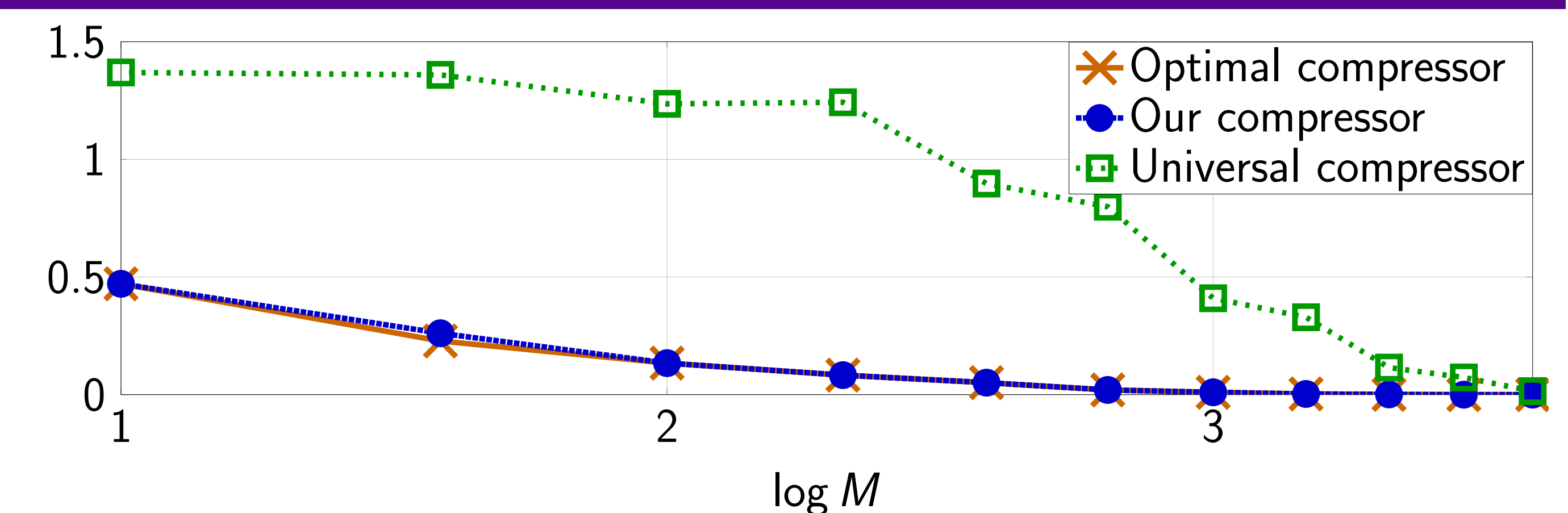
- ▶ optimal compressor f^* — when feasible to compute, i.e, small $|\mathcal{X}|$;
- ▶ our KL-greedy compressor;
- ▶ universal compressor from [2], which is designed for reconstruction under log-loss distortion.

For the empirical type-II error rate, consider a threshold T such that type-I error rate $< \epsilon = 0.05$ for a given compressor at rate M .

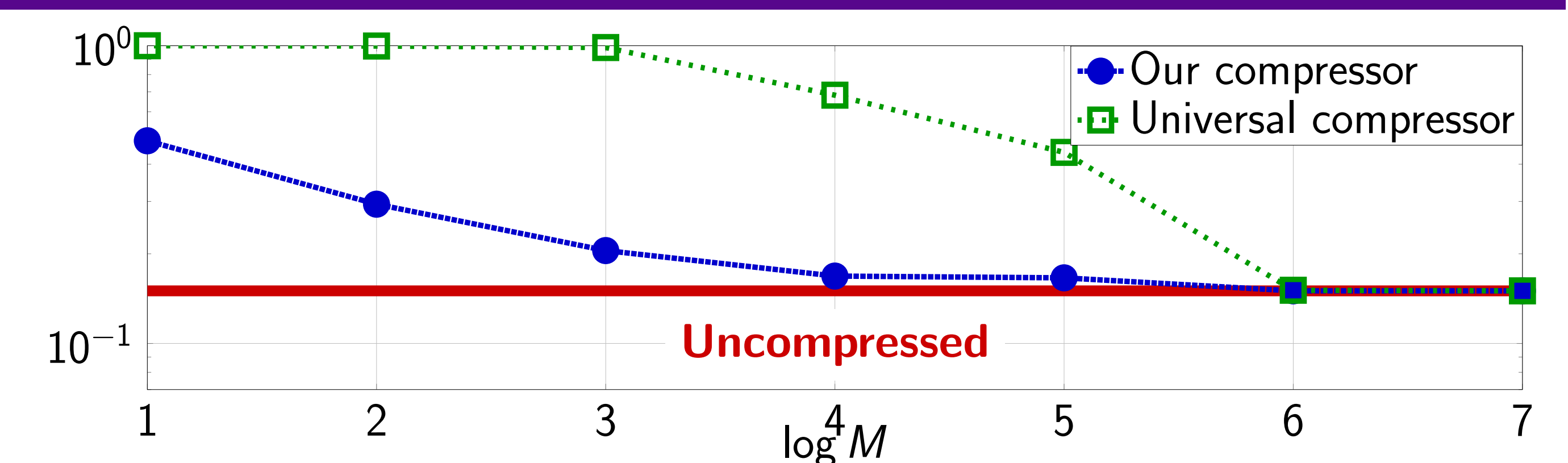
Results: Distributions and Compressor for $|\mathcal{X}| = 13, M = 4$



Results: Compression penalty $\Delta_f(R)$ for $|\mathcal{X}| = 13$



Results: Type-II Error Rate for $|\mathcal{X}| = 256, n = 5, \epsilon = 0.05$



Conclusions

- ▶ Formulation for the optimal compressor for hypothesis testing.
- ▶ Proposed the empirical “KL-greedy” compressor: it can be computed in polynomial time and preserves the *useful* information.
- ▶ Task-aware compression achieves error rate comparable to the uncompressed case for low rates.

References

- [1] T. M. Cover and J. A. Thomas, *Elements of Information Theory (Wiley Series in Telecommunications and Signal Processing)*. USA: Wiley-Interscience, 2006.
- [2] Y. Shkel, M. Raginsky, and S. Verdú, “Universal lossy compression under logarithmic loss,” in *2017 IEEE International Symposium on Information Theory (ISIT)*, Jun. 2017, pp. 1157–1161.
- [3] F. Carpi, S. Garg, and E. Erkip, “Single-shot compression for hypothesis testing,” in *22nd IEEE Int. Workshop on Signal Processing Advances In Wireless Communications (SPAWC)*, Sep. 2021.

Acknowledgements

This work was supported in part by NSF-Intel grant #2003182 and NSF grant #1925079.