

on
(mostly lossless)
compression induced learning

Israel Statistics and Data Science Annual Conference

June 2026

Tsachy Weissman
Stanford University

my work

- information theory
- information processing
- data compression

what this talk is not about

- compression of neural networks
- compression of gradients while training
- KV cache compression
- etc.

outline

- apology
- lossless compression landscape
- compressors as Sequential Probability Assignments (SPAs)
- the LZ78 SPA: classical and modern twists
- example applications: classification and “Gen AI”
- LZ78 as a data generator
- the interplay between lossy compression and denoising
- textual transform coding
- some takeaways

(will try to accommodate a diversity of backgrounds)

apology in advance

there will be no justice to the (classical and recent) literature but it is served in:

- [SW26]
- [Omri et al. 25]
- [Gorle et al. 26]
- [SDHW26]
- [SOW26]
- [Ding et al. 26]
- [W23]

lossless compression landscape in practice

Large Text Compression Benchmark

Matt Mahoney
Last update: Oct. 8, 2024. history

This competition ranks lossless data compression programs by the compressed size (including the size of the decompression program) of the first 10^9 bytes of the XML text dump of the English version of Wikipedia on Mar. 3, 2006. About the test data.

The goal of this benchmark is not to find the best overall compression program, but to encourage research in artificial intelligence and natural language processing (NLP). A fundamental problem in both NLP and text compression is modeling: the ability to distinguish between high probability strings like recognize speech and low probability strings like reckon eyes peach. Rationale.

This is an open benchmark. Anyone may contribute results. Please read the rules first.

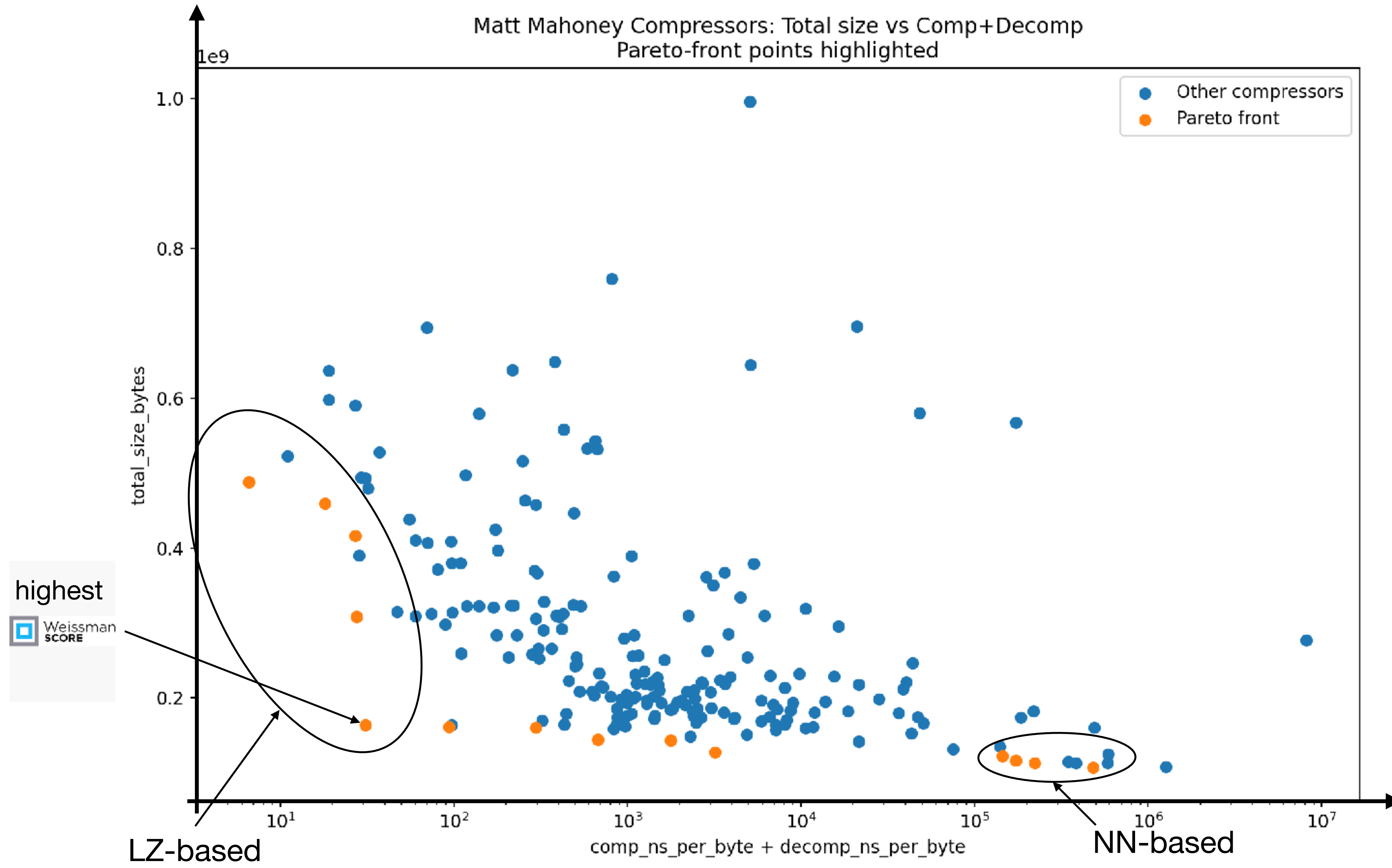
Open source compression improvements to this benchmark with certain hardware restrictions may be eligible for the Humer Prize.

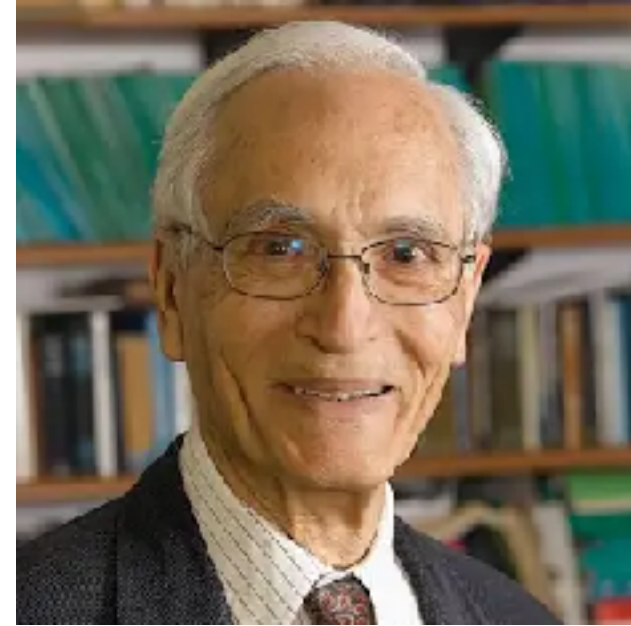
Benchmark Results

Compressors are ranked by the compressed size of enwik9 (10^9 bytes) plus the size of a zip archive containing the decompressor. Options are selected for maximum compression at the cost of speed and memory. Other data in the table does not affect rankings. This benchmark is for informational purposes only. There is no prize money for a top ranking. Notes about the table:

- Program: The version believed to give the best compression. A | denotes a combination of 2 programs.
Compression options: selected for what I believe gives the best compression.
enwik8: compressed size of first 10^8 bytes of enwik9. This data is used for the Humer Prize, and is also ranked here but has no effect on this ranking.
enwik9: compressed size of first 10^9 bytes of enwik9-20060303-pages-articles.xml.
decompressor size: size of a zip archive containing the decompression program (source code or executable) and all associated files needed to run it (e.g. dictionaries). A letter following the size has the following meaning:
x = executable size.
s = source code size (if available and smaller).
d = size of a separate decompression program (separate from compression). For self extracting archives (SFX), the size is 0 because the decompressor and compressed data are combined into one file.
For testing, if no zip file is supplied I create archives using InfoZIP 2.32-9. (Prior to July 1, 2008 I used 7zip 4.32 -trip -mx-9).
Total size: total size of compressed enwik9 + decompressor size, ranked smallest to largest.
Comp: compression rate in nanoseconds per byte on the largest file tested (e.g. seconds for enwik9). Speed is approximate and has no effect on ranking. A ~ means 'very approximate'. Not all tests are done on the same computer. Times reported are the smaller of process time (summed over processors if multi-threaded) or real time as measured with timer. If there is no note then the program was tested on a Compaq Presario 1440, 2.108 GHz, Amlion-64 3500+ in 32 bit Windows XP. An underlined time means that no better compressor is faster.
Decomp: decompression time as above. If blank, decompression was not tested yet and ranking is pending verification that the output is identical. An underlined time means that no better compressor is faster.
Mem: approximate memory used for compression in MB. Decompression uses the same or possibly less. There is some ambiguity whether a megabyte means 10^6 bytes or 2^20 bytes. The approximation is coarse enough that it doesn't matter. I use peak memory as measured with Windows Task Manager during compression (so if you really want to know, 1 MB = 1,024,000 bytes.) Memory does not include swap or temporary files. An underlined value means that no better compressor uses less memory.
Alg: compression algorithm, referring to the method of parsing the input into symbols (strings, bytes, or bits) and estimating their probabilities (modeling) for choosing code lengths. Symbols may be arithmetic coded (fractional bit length for best compression), Huffman coded (bit aligned for speed), or byte aligned as a preprocessing step.
Dict (Dictionary). Symbols are words, coded as 1 or 2 bytes, usually as a preprocessing step.
LZ (Lempel Ziv). Symbols are strings.
LZ77: repeated strings are coded by offset and length of previous occurrence.
LZW (LZ Welch): repeats are coded as indexes into a dynamically built dictionary.
ROLZ (Reduced Offset LZ): LZW with multiple small dictionaries selected by context.
LZP (LZ predictive): ROLZ with a dictionary size of 1.
on (Order-n, e.g. o0, o1, o2...): symbols are bytes, modeled by frequency distribution in context of last n bytes.
PPM (Prediction by Partial Match): order-n, modeled in longest context matched, but dropping to lower orders for byte counts of 0.
SR (Symbol Ranking): order-n, modeled by time since last seen.
HWT (Huffman Wheel Transform): bytes are sorted by context, then modeled by order-0 SR.
ST (Sort Transform): HWT using stable sort with truncated string compression.

Table with columns: Program, Compression Options, Compressed size (enwik8, enwik9), Decompressor size (zip), Total size (enwik9+prog), Time (ns/byte) Comp, Decomp, Mem, Alg, and Rank. It lists various compression programs like nroep v3.7, xz v2.1, 7z2-mkx, etc., with their respective performance metrics.





Jacob Ziv

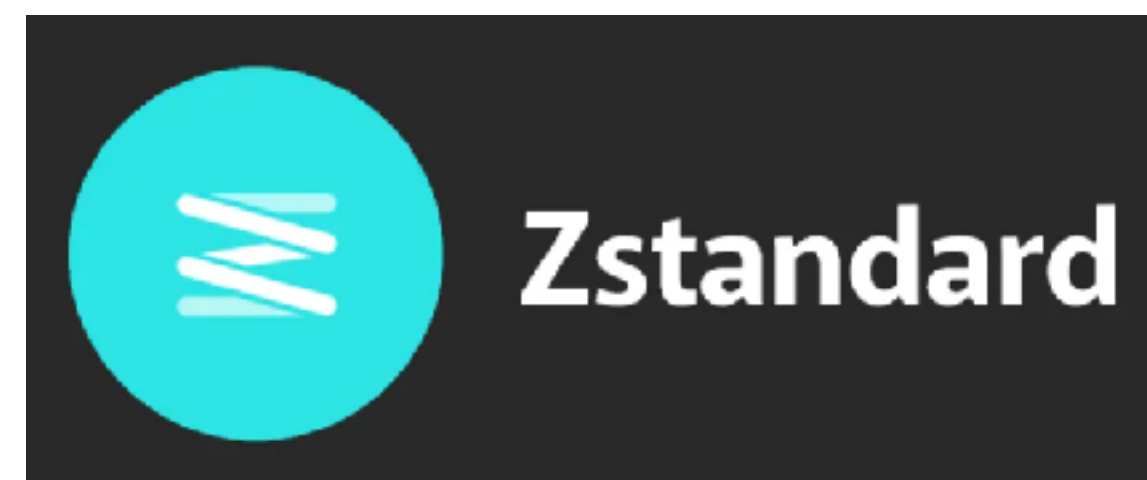


Abraham Lempel

- LZ77 (sliding window)
- LZ78 (incremental parsing)

LZ compression achieves:

- the fundamental limits (entropy rate and finite-state compressibility)
- (near) linear complexity
- elegance
- state of the art performance on real data

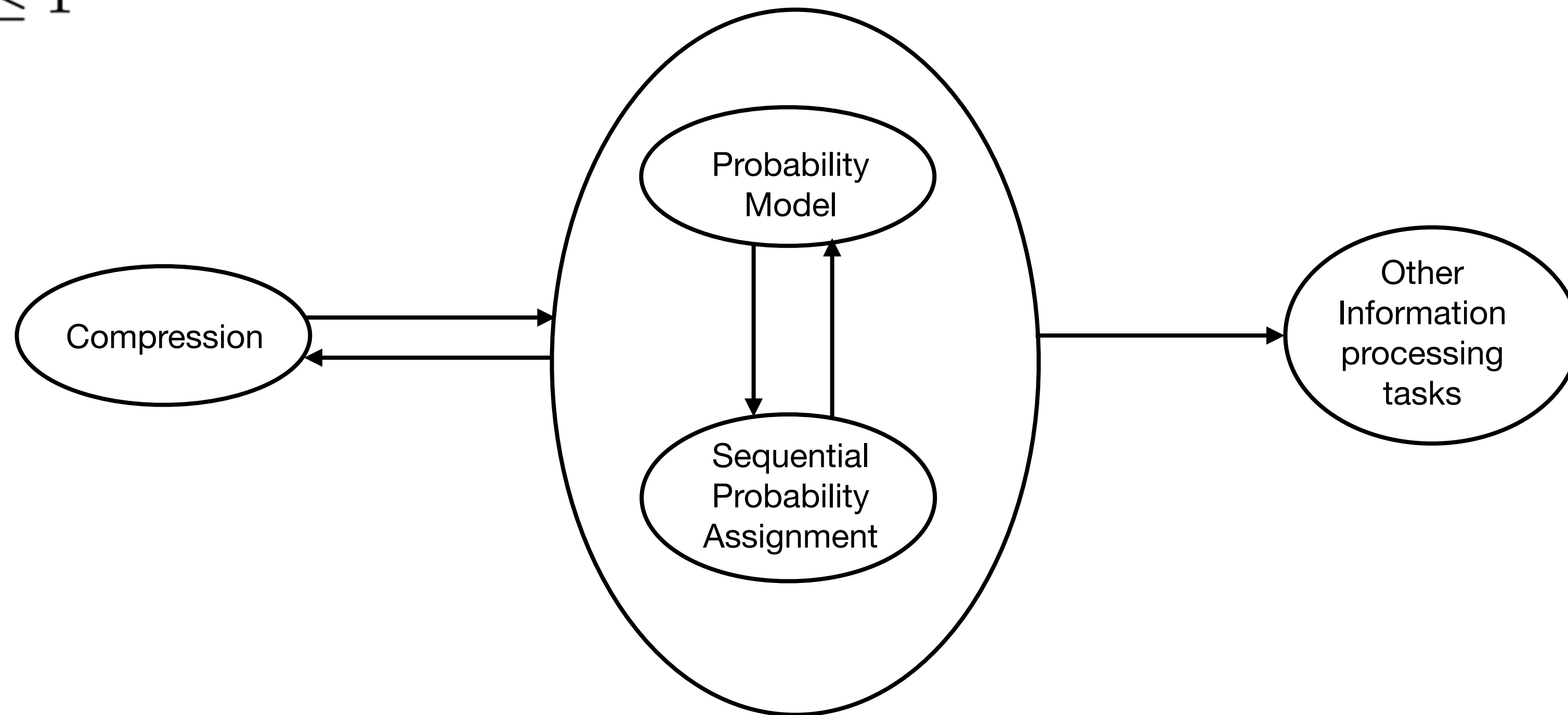


why care about (lossless) compression?

in(formation) theory (I0I)

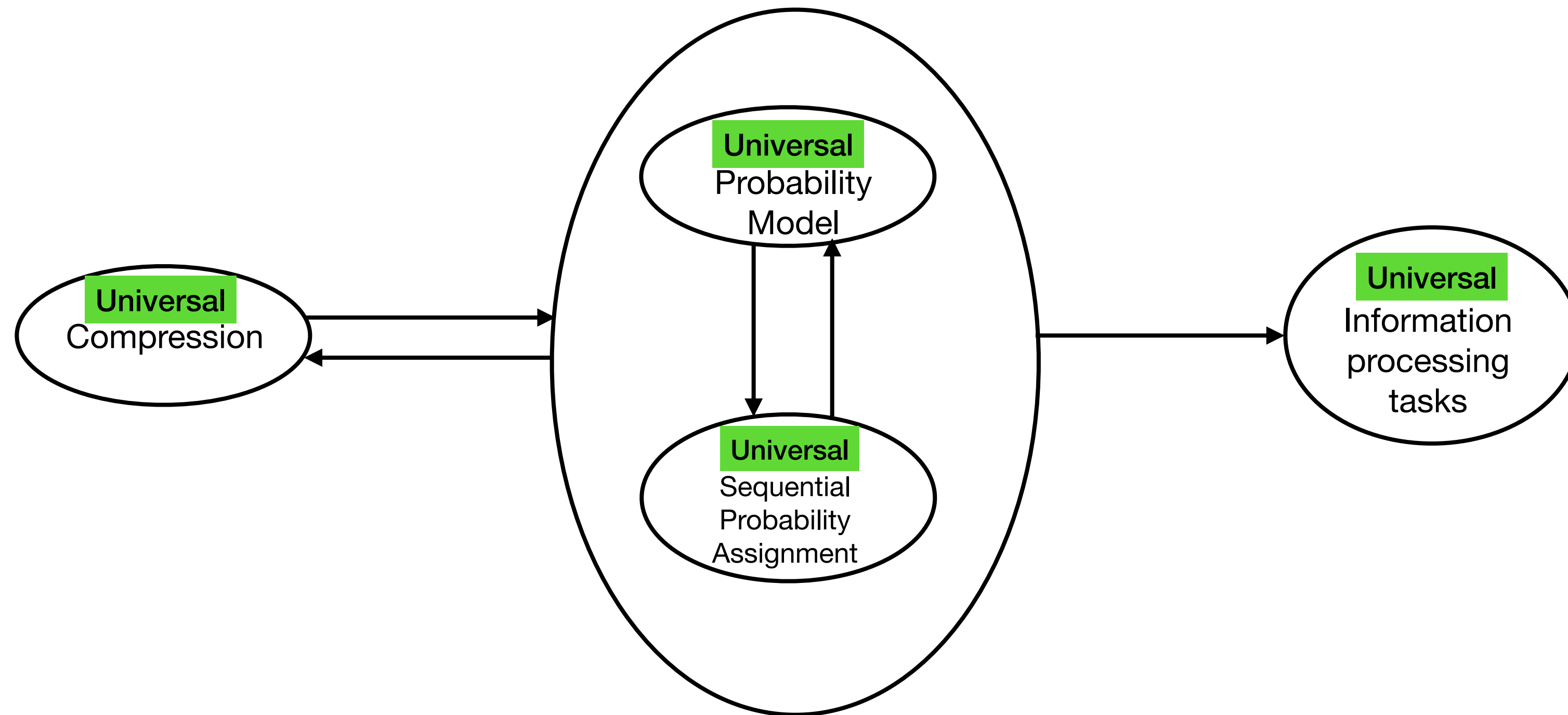
Kraft inequality:

$$\sum_{\text{data}} 2^{-\ell(\text{data})} \leq 1$$

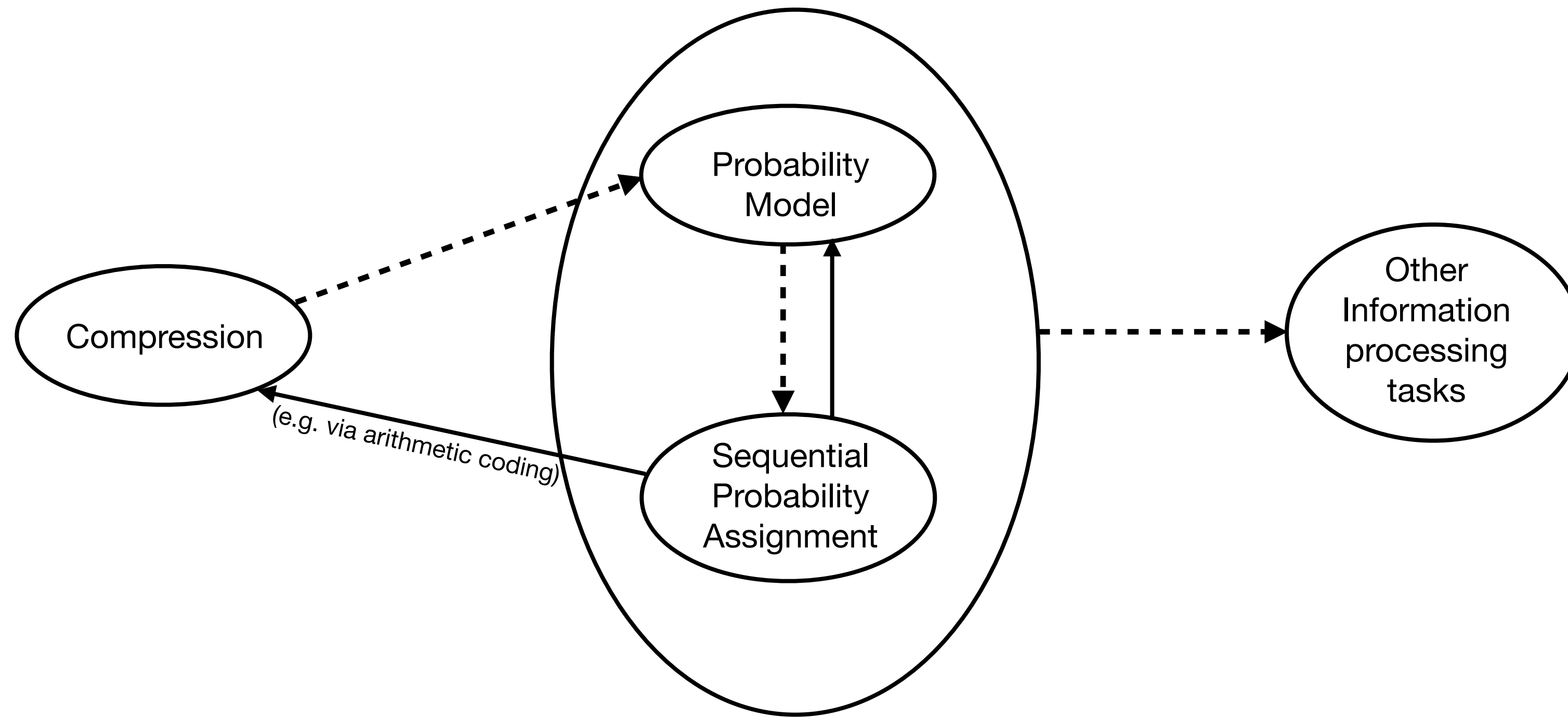


- (gets rediscovered periodically in other communities, cf., e.g., “Language modeling is compression”, ICLR 2024)

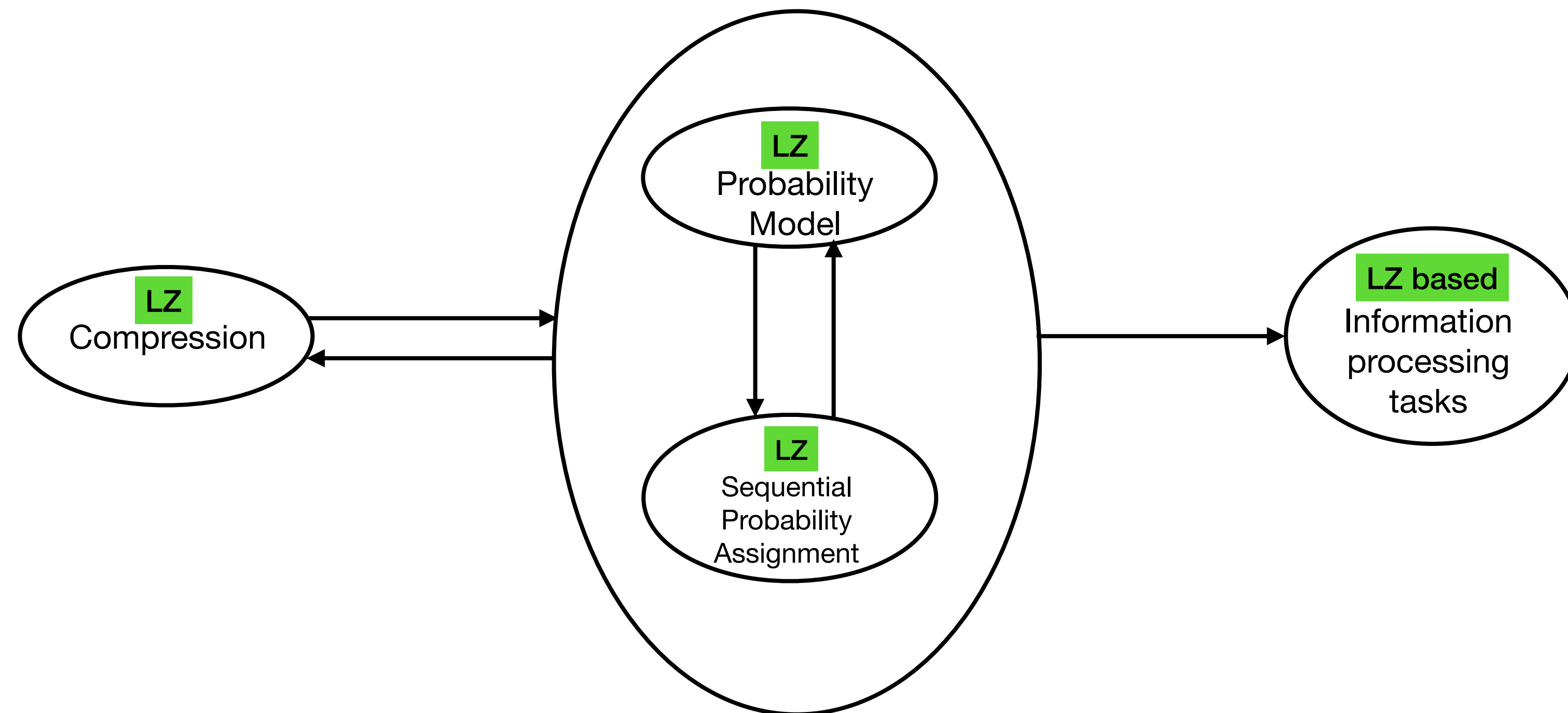
universality is also inherited



practically/algorithmically



LZ(78): both universality and practicality is inherited

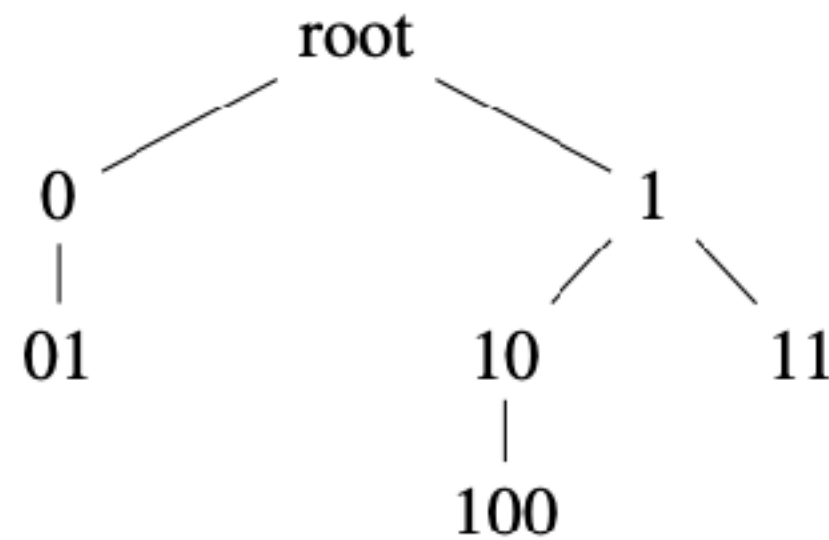
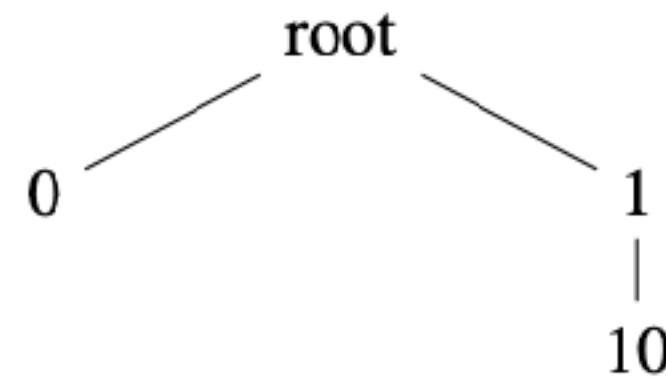
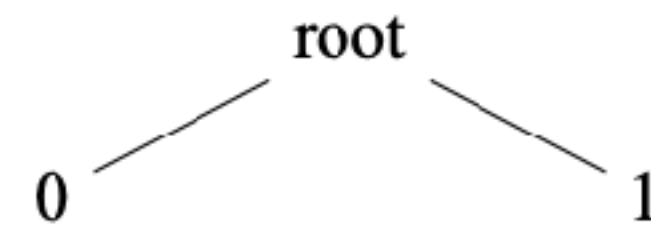


LZ78 induced Sequential Probability Assignment (SPA)

$x^n = 01100110011$

incremental parsing: $0, 1, 10, 01, 100, 11$

induces tree:



tree induces SPA

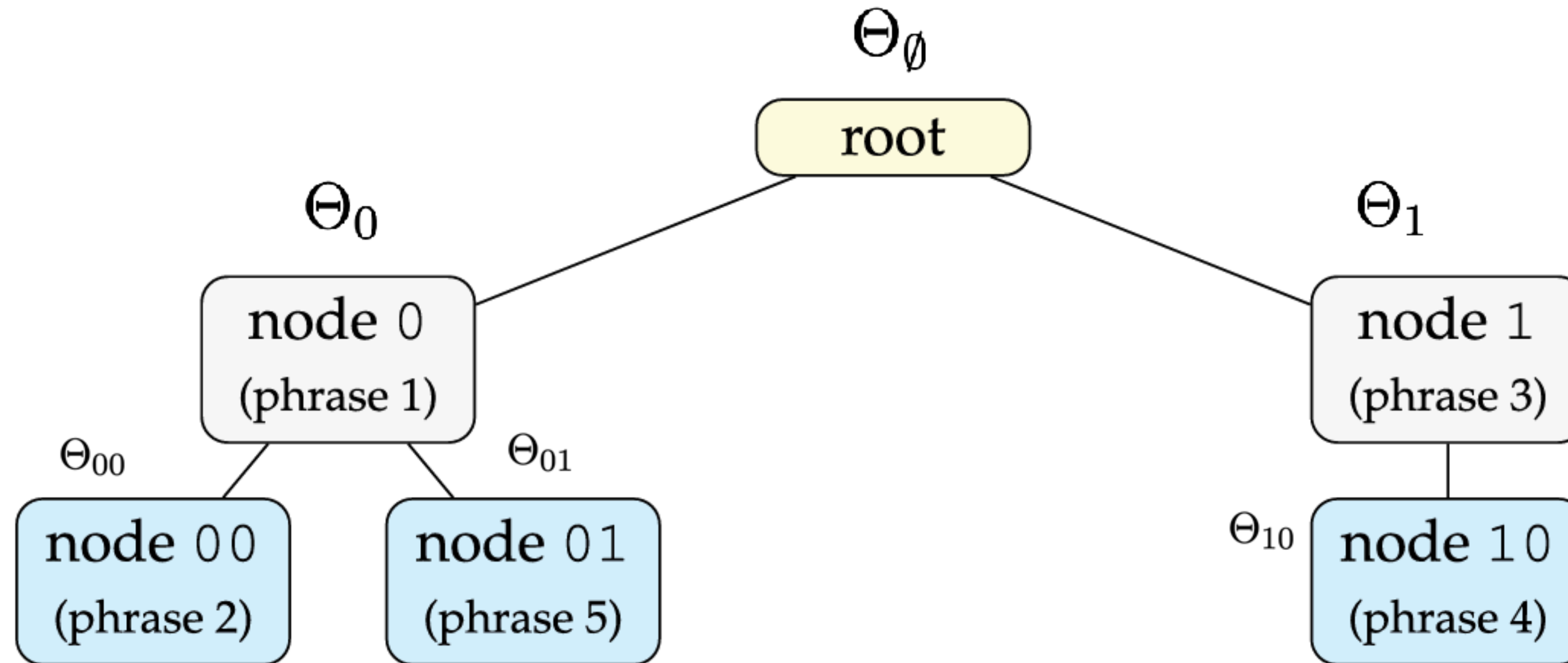
of phrases in the parsing

$$\ell_{LZ}(x^n) = \frac{C(x^n) \log C(x^n)}{n} + o(1)$$

N. Sagan and T. Weissman,
"A Family of LZ78-based Universal Sequential Probability Assignments",
Information Theory transactions, April 2026
[SW26]

00011001 \rightarrow 0,00,1,10,01,...

$$\Theta. \sim f$$



example

when f is the Dirichlet prior:

$$q^f(x_t | x^{t-1}) = \frac{N(x_t | x^{t-1}) + \gamma}{(t-1) + \gamma A}$$

$$q^{LZ78, f}(a | x^{t-1}) = \frac{N(a | x^{t-1}, z_c(x^{t-1})) + \gamma}{\sum_{b \in \mathcal{A}} N(b | x^{t-1}, z_c(x^{t-1})) + \gamma A}$$

$z_c(x^{t-1})$ is the node in the LZ tree reached at time $t-1$

(aka “finite-state predictability under log loss”)

Definition IV.1 (Finite-State SPA). A finite-state SPA is such that the probabilities assigned depend solely on the current state of an underlying finite-state mechanism. Concretely, q is a M -state SPA if \exists next-state function $g : [M] \times \mathcal{A} \rightarrow [M]$, prediction function $f : [M] \rightarrow \mathcal{M}(\mathcal{A})$, and initial state $s_1 \in [M]$ such that

$$\forall t \geq 1, q(\cdot|x^{t-1}) = f(s_t), \quad \text{and} \quad s_{t+1} = g(s_t, x_t).$$

The set of all M -state SPAs is denoted \mathcal{F}_M .

The optimal log loss of any finite-state SPA is defined as follows:

Definition IV.2 (Optimal Finite-State Log Loss). The minimum log loss of a M -state SPA for sequence \mathbf{x} is defined as

$$\lambda_M(\mathbf{x}) \triangleq \overline{\lim}_{n \rightarrow \infty} \lambda_M(x^n) \quad \text{for} \quad \lambda_M(x^n) \triangleq \min_{q \in \mathcal{F}_M} \frac{1}{n} \sum_{t=1}^n \log \frac{1}{q(x_t|x^{t-1})}.$$

The optimal finite-state log loss takes the number of states to infinity,

$$\lambda(\mathbf{x}) \triangleq \lim_{M \rightarrow \infty} \overline{\lim}_{n \rightarrow \infty} \lambda_M(x^n).$$

$\lambda_M(\mathbf{x})$ is monotonically non-increasing in M and bounded below, so the outer limit is guaranteed to exist.

(aka “Markov predictability under log loss”)

The definition of the optimal Markov SPA log loss is analogous to that of the optimal finite-state SPA log loss:

Definition IV.6 (Optimal Markov Log Loss). For any sequence \mathbf{x} , the optimal log loss of a k -order Markov SPA is

$$\mu_k(\mathbf{x}) \triangleq \overline{\lim}_{n \rightarrow \infty} \mu_k(x^n) \quad \text{for} \quad \mu_k(x^n) \triangleq \min_{q \in \mathcal{M}_k} \frac{1}{n} \sum_{t=1}^n \log \frac{1}{q(x_t | x^{t-1})}.$$

The optimal Markov SPA log loss is defined by taking the context length to ∞ :

$$\mu(\mathbf{x}) \triangleq \lim_{k \rightarrow \infty} \mu_k(\mathbf{x}) = \lim_{k \rightarrow \infty} \overline{\lim}_{n \rightarrow \infty} \mu_k(x^n).$$

As with $\lambda(\mathbf{x})$, the outer limit is guaranteed to exist.

Theorem: $\lambda(\mathbf{x}) = \mu(\mathbf{x})$

trivial

surprising

$$\lambda(\mathbf{x}) \leq \mu(\mathbf{x})$$

$$\lambda(\mathbf{x}) \geq \mu(\mathbf{x})$$

universality of the LZ78 family of SPAs

Theorem V.1 (Universality of LZ78 SPA). *For prior distribution Π such that $\text{supp}(\Pi) = \mathcal{M}(\mathcal{A})$, $q^{\text{LZ78}, \Pi}$ from Construction III.5 satisfies, for any individual sequence,*

$$\overline{\lim}_{n \rightarrow \infty} \frac{1}{n} \log \frac{1}{q^{\text{LZ78}, \Pi}(x^n)} \leq \lambda(\mathbf{x}).$$

universality in the stochastic setting (classical)

entropy rate: $\mathbb{H}(\mathbf{X}) = \lim_{n \rightarrow \infty} \frac{1}{n} H(X^n)$

Theorem E.2. *Suppose \mathbf{X} is a stationary process and the SPA q satisfies*

$$\overline{\lim}_{n \rightarrow \infty} \frac{1}{n} \log \frac{1}{q(x^n)} \leq \mu(\mathbf{x}),$$

for all individual sequences \mathbf{x} . Then, if $\mathbb{H}(\mathbf{X})$ is the entropy rate of the stochastic process,

$$\lim_{n \rightarrow \infty} \mathbb{E} \left[\frac{1}{n} \log \frac{1}{q(X^n)} \right] = \mathbb{H}(\mathbf{X}).$$

(equivalently: $\lim_{n \rightarrow \infty} \frac{1}{n} D(\text{True law of } X^n \parallel \text{law of } X^n \text{ under } q) = 0$)

When \mathbf{X} is also ergodic: $\mu(\mathbf{X}) = \lambda(\mathbf{X}) = \mathbb{H}(\mathbf{X}) \quad a.s.$

result at the heart

Theorem III.11. *For any prior such that $\text{supp}(\Pi) = \mathcal{M}(\mathcal{A})$,*

$$\lim_{n \rightarrow \infty} \max_{x^n} \left| \frac{1}{n} \log \frac{1}{q^{\text{LZ78}, \Pi}(x^n)} - \frac{C(x^n) \log C(x^n)}{n} \right| = 0.$$

where $C(x^n)$ is the number of phrases in the LZ parsing of x^n and

$$\ell_{\text{LZ}}(x^n) = \frac{C(x^n) \log C(x^n)}{n} + o(1)$$

initial experiments in [SW26]

Table I
RESULTS OF CLASSIFICATION EXPERIMENTS USING THE LZ78 SPA.

Dataset	MNIST	Fashion-MNIST	IMBD	Spam Emails
Accuracy (%)	75.36	72.16	75.62	98.12
Training Time (s)	14	15	16	14

initial experiments in [SW26] (cont.)

This is the moon with the fairest charge thee stay;
Which is not so.

Provost:
Art thou not
That which withal; you go you to Baptista; or, but I am not Lucentio,
Red the beasts, that Warwick,
And those that runatice of his auture of our straight and will play the his profane eyes, came to see thy servant
so die.

LUCIO:
But what lives in Signior Gremio: fools
His will I may have auture of his auture of our common good time
Unto the Tower,
Give me thy hand as come enough
And then he shall not rather with the Lord Stand be so longer see that are fond, as thou say the orace to speak
brother, or oints?

BIONDELLO:
I may more straight and will play the his punish his convey much this leavine art thou that will not should be
thus sir, there's face.
Go you to Baptista; or, lo, here all abroad in the

Figure 1. 800 symbols generated via the LZ78 SPA, trained on the `tiny-shakespeare` dataset, seeded with the sequence "This."

Y. Omri, N. Sagan, E. Min, H. Choi, T. Moon and T. Weissman,
"Genomic Data Classification via Universal Compression",
[Omri et al. 25]
ISMB 2025 and
submitted for journal publication
<https://doi.org/10.21203/rs.3.rs-6363017/v1>

main idea

- compress the labeled data from each class jointly with the test sequence
- choose class with which the test sequence jointly compresses best
- when using LZ78:
 - training is linear in the training data
 - inference is linear in the test data
 - induced LZ78 tree gives a succinct and transparent summary of the class statistics

on accuracy:

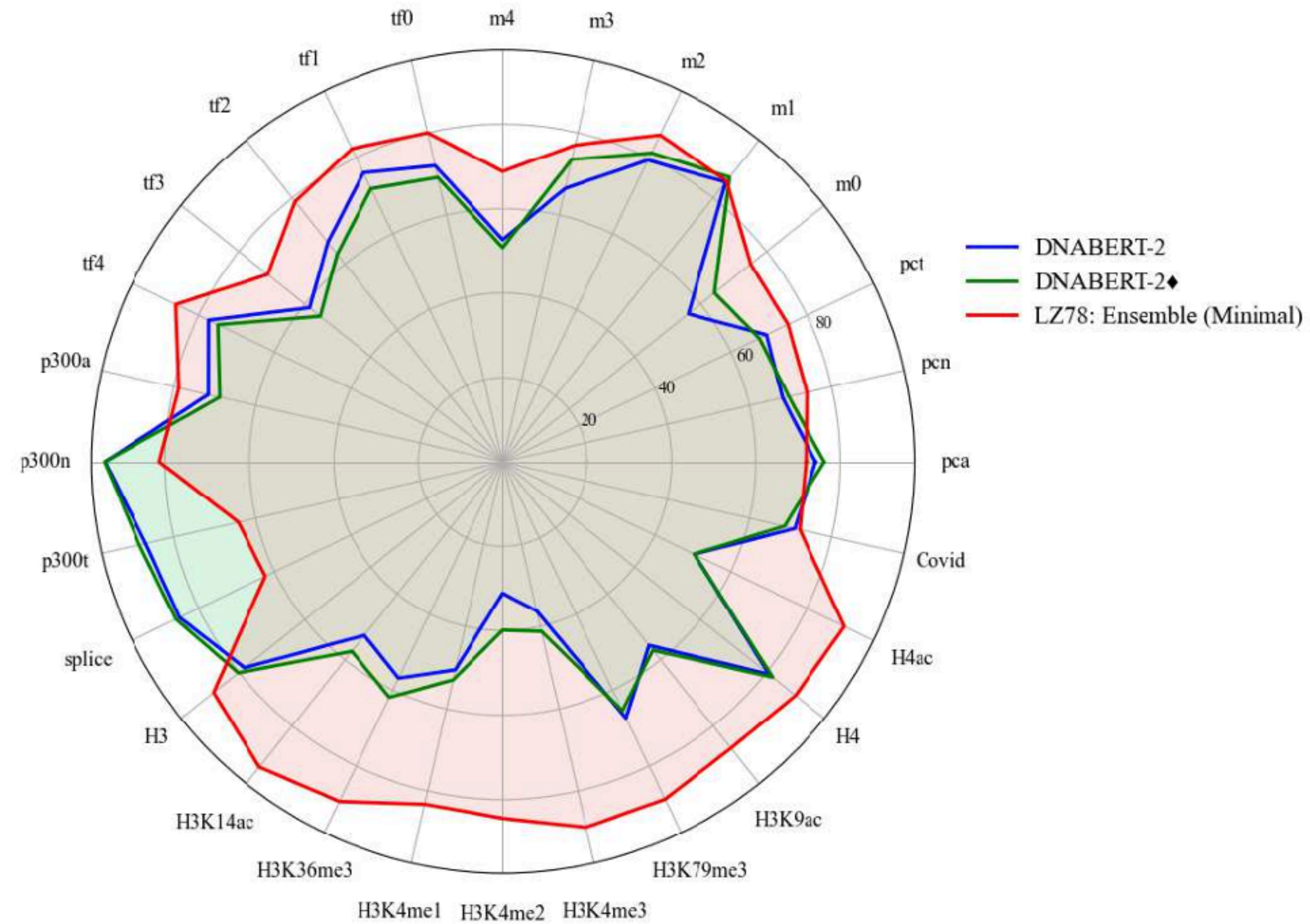
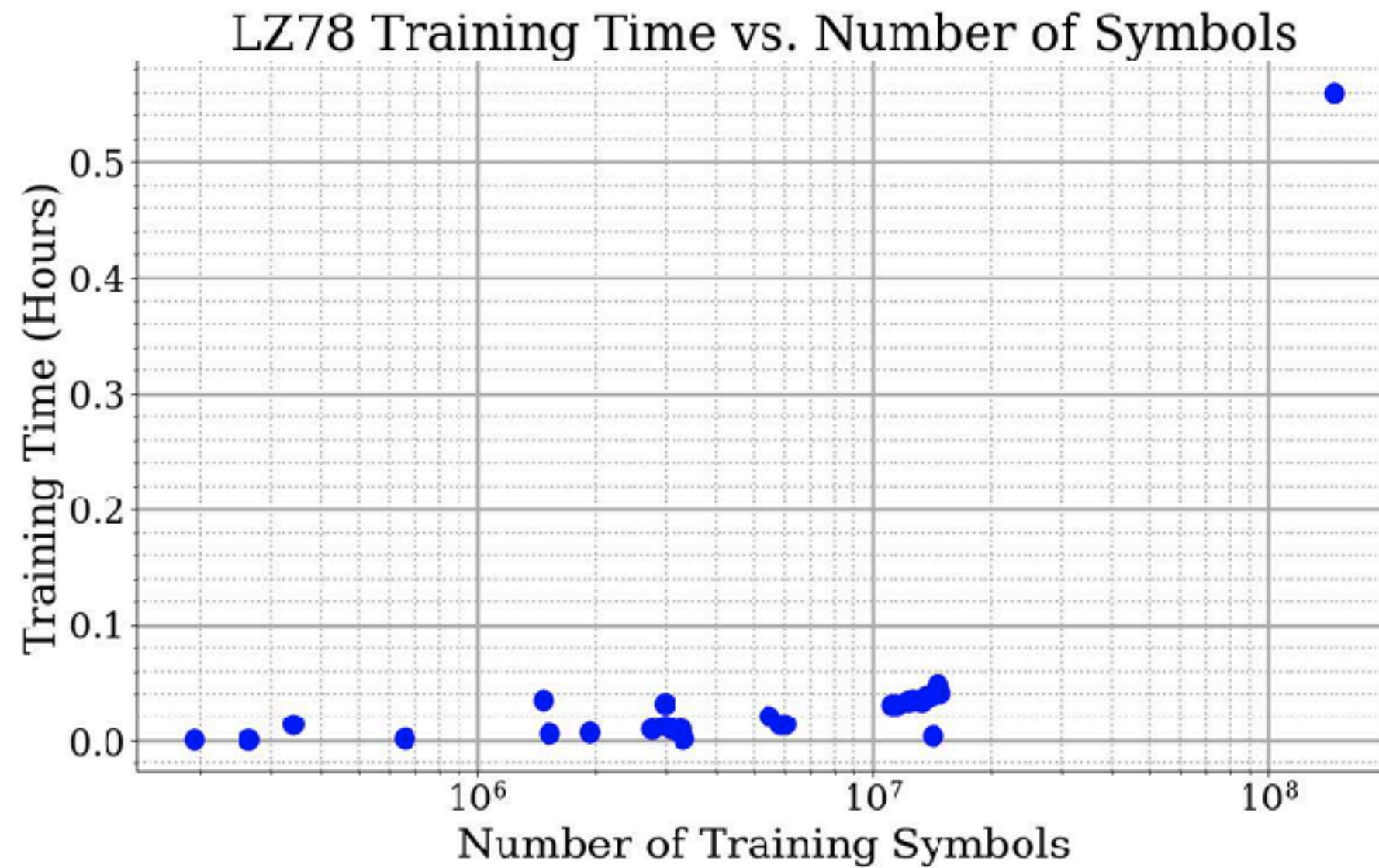


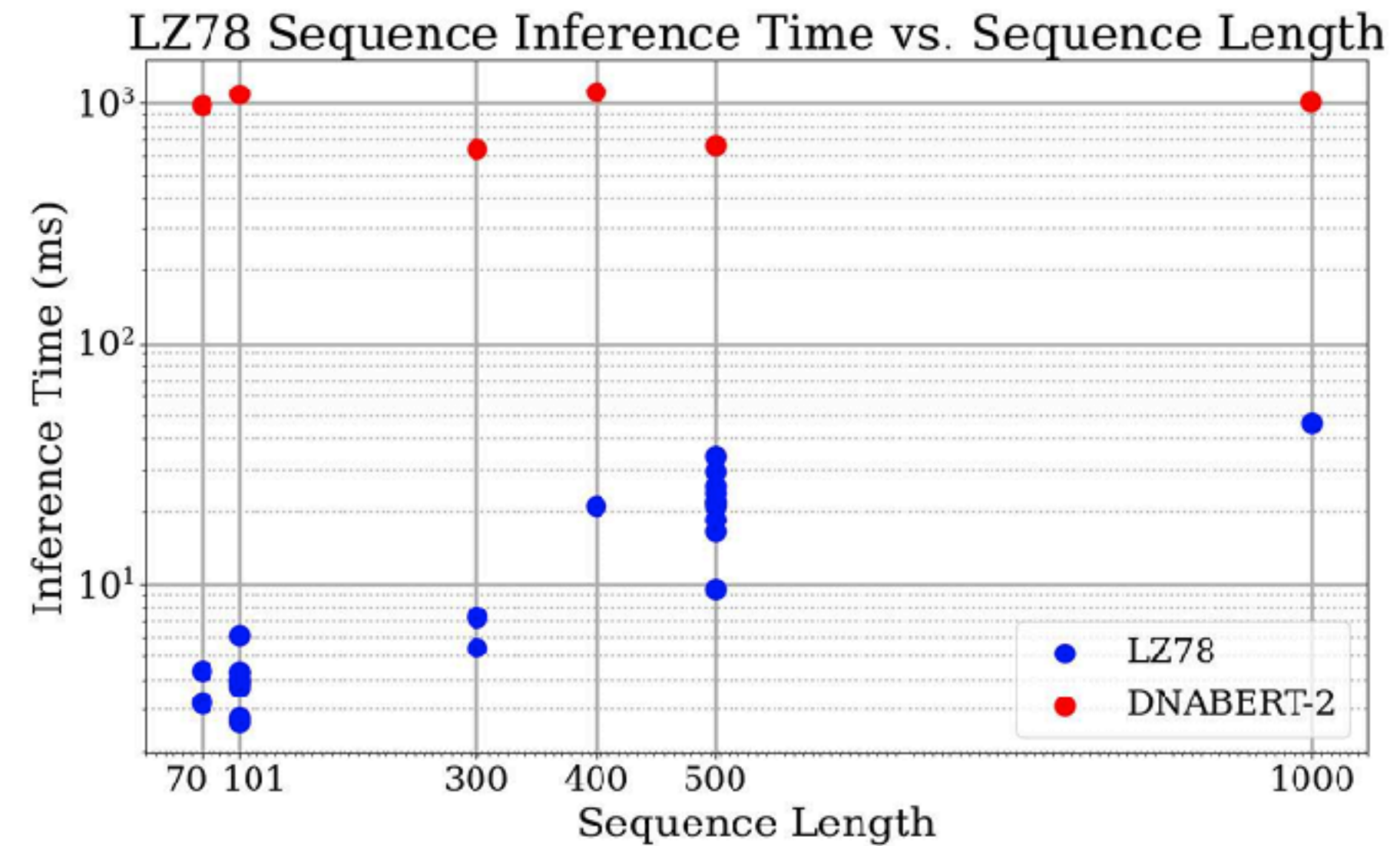
Fig. 2: Accuracy results on the GUE benchmark dataset

(GUE: Genomic Understanding Evaluation Suite)

on compute:



(a) Efficiency of LZ78 Classifier Training; average training time of ≈ 2.3 minutes per dataset on a modern CPU.



(b) Inference time speedup of Using an LZ78-Based classifier compared to DNABERT-2.

merits

- improved accuracy
- dramatically faster
- requires substantially less training data
- transparent white-box logic
- theoretically justified

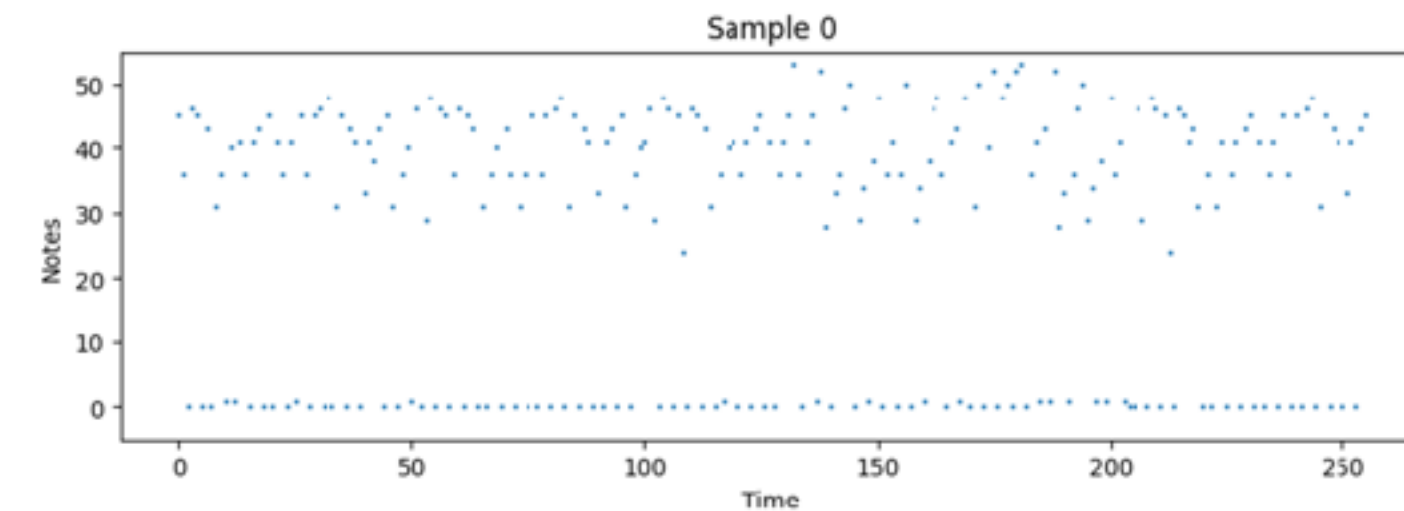
A. Gorle, C. Ding, S. Bhattacharya, D. Haster, N. Sagan and T. Weissman,
"LZMidi: Compression-Based Symbolic Music Generation",
ICASSP 2026
[Gorle 26]

LZMidi

- **Goal:** Offer an alternative to large ML-based generators with a *universal, CPU-only* compression-driven generator (sampler)
- **Core idea:** Fast sampling from the sequential probability assignment induced from an LZ78 tree (compression)
- **Application:** Symbolic Music generation
- **Theoretical results:** Convergence to true source statistics

Training and Dataset:

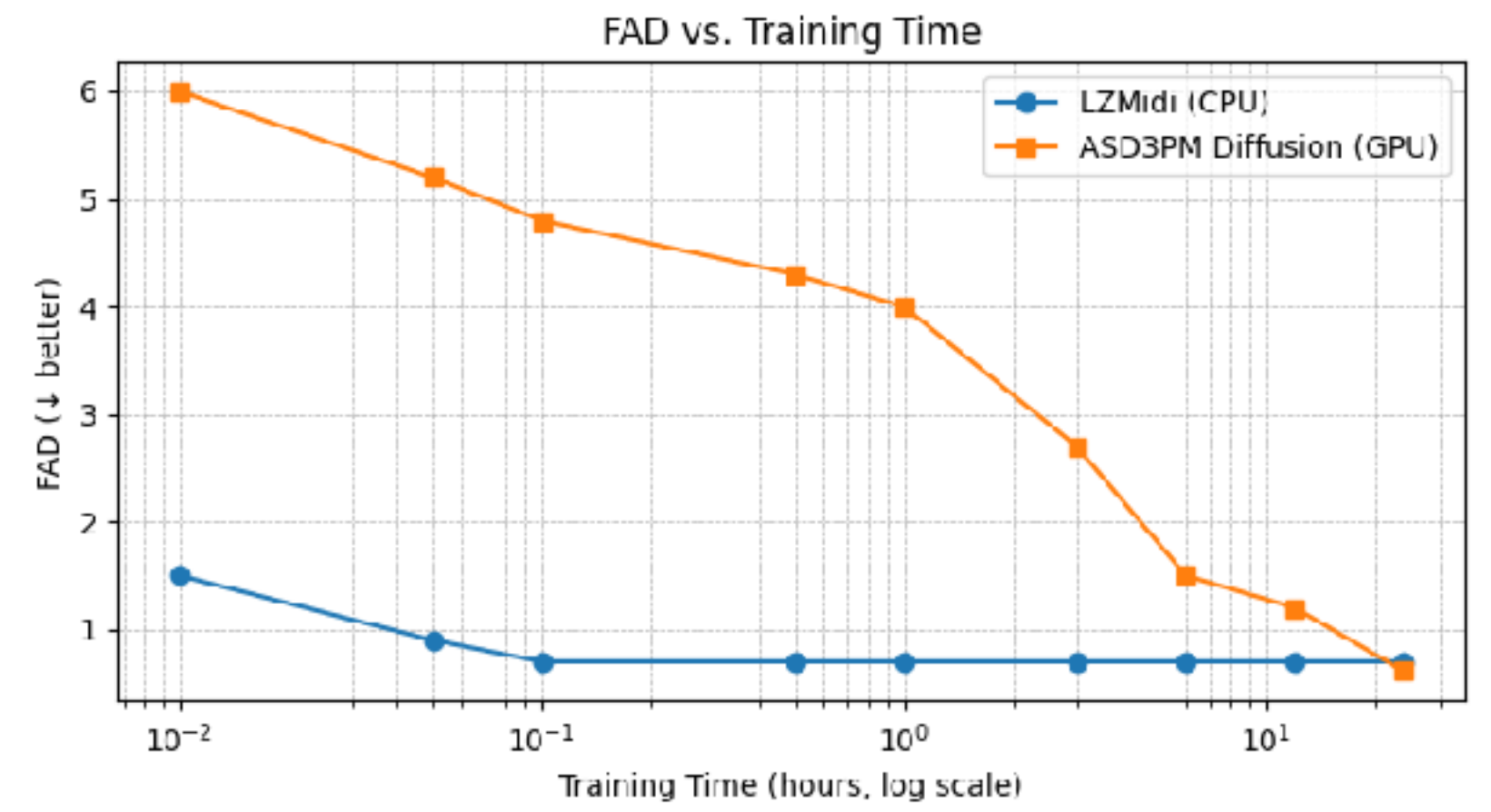
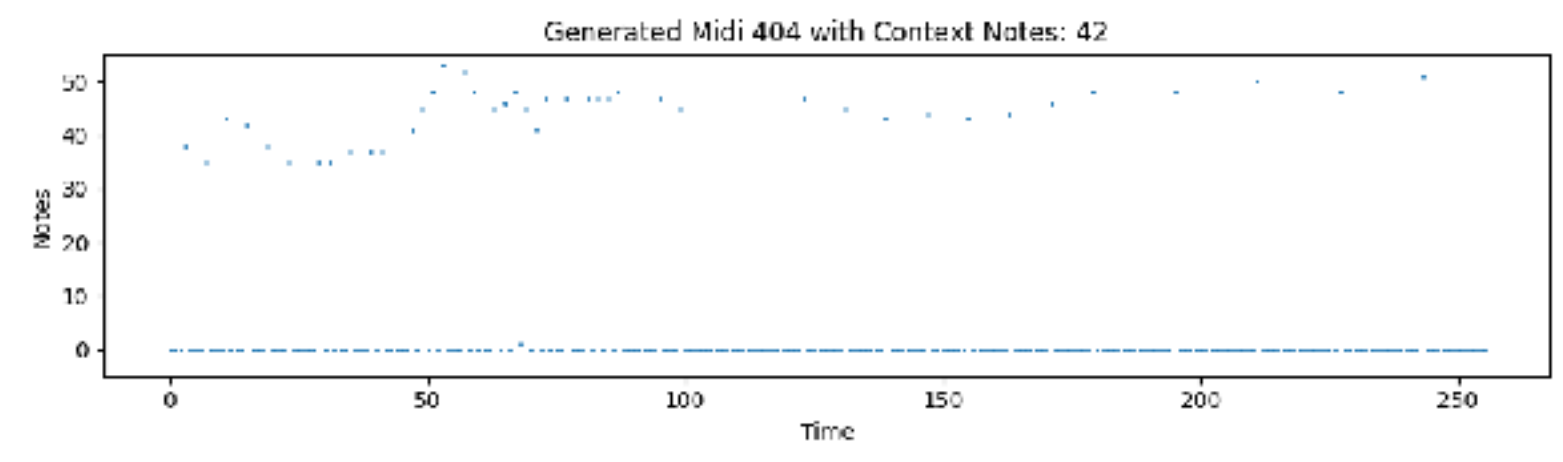
- Lakh MIDI Dataset (LMD): 648,574 symbolic MIDI sequences (each 256 tokens over alphabet {0-89}, where 0=rest, 1=continuation, 2-89=pitches)
- 80/20 split: ~518 k sequences for training (used to build/update the LZ78 tree)
- Training: a simple one-pass through all training sequences to build the LZ tree, CPU-friendly



LZMidi highlights

- **Sampling/generation:** initialized by a single randomly selected symbol (a direct child of the LZ tree root) from the 90-symbol alphabet
- Reference sample on the right
- **Metrics (vs. deep models):**
 - FAD as low as $\downarrow 0.64$ (comparable to MusicVAE and fully-trained Diffusion baselines)
 - Training: 30 \times faster, Generation: 300 \times faster (CPU only) versus neural baselines (GPU)
 - substantially outperforms compute-matched deep-learning baselines
- **Compute Efficiency:** Hits FAD ~ 0.7 within minutes on M1 CPU vs hours of GPU training for deep-learning models

	Train (s)	Generation (s/sample)	Model Size, MB	TE (kJ)	GE (J/sample)
LZMidi	107.7	0.016	287.1	9.144	1.36
ASD3PM	3480	5.4	306.2	2088	3240



Key Takeaway: Quality comparable to SOTA, but *minutes on CPU* instead of *GPU-days*

examples:



got us thinking

- about compressors and SPAs as information sources
- for some it's very clear: Markov, CTW,...
- what about the LZ78 SPA as an information source?

N. Sagan, A. Dembo, M. Ho and T. Weissman,
"The LZ78 Source", Information theory transactions, 2026 [SDHW26]

question:

what are the distributional and entropic properties of a process generated from the LZ78 source?

[SDHW26]:

Theorem III.11 (Shannon-McMillan-Breiman-Type Result). *Let $\mathbf{X} \sim Q^{LZ, \Pi}$ be generated from the LZ78 source with $\text{supp}(\Pi) = \mathcal{M}(\mathcal{A})$. Then, almost surely,*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \frac{1}{Q_{X^n}^{LZ, \Pi}(X^n)} = \mathbb{E}_{\Theta \sim \Pi} [H(\Theta)].$$

Theorem III.12 (Entropy Rate). *The entropy rate of the LZ78 source, under any Π , is*

$$\lim_{n \rightarrow \infty} \frac{1}{n} H(X^n) = \mathbb{E}_{\Theta \sim \Pi} [H(\Theta)].$$

on the other hand:

Theorem III.13 (Finite-state Predictability). *For \mathbf{X} generated from the LZ78 source, the finite-context Markov predictability is almost surely*

$$\mu(\mathbf{X}) = H(\mathbb{E}_{\Theta \sim \Pi} [\Theta]).$$

so there's a "Jensen gap" (in stark contrast to stationary ergodic processes)

[SDHW26]:

divergence properties

When Π has full support, among the results of [6] is the universality of $Q^{\text{LZ},\Pi}$ in the sense that

$$\lim_{n \rightarrow \infty} \frac{1}{n} D \left(P_{X^n} \parallel Q_{X^n}^{\text{LZ},\Pi} \right) = 0 \quad \forall P \text{ stationary.}$$

In contrast:

Theorem III.6 (Divergence with Markovian Law: Pointwise Result). *Let \mathbf{X} be generated from the LZ78 source $Q^{\text{LZ},\Pi}$ with $\text{supp}(\Pi) = \mathcal{M}(\mathcal{A})$. For any k^{th} -order Markov law $P \in \mathcal{M}_k$, almost surely*

$$\lim_{n \rightarrow \infty} \frac{1}{n} \log \frac{Q_{X^n}^{\text{LZ},\Pi}(X^n)}{P_{X^n}(X^n)} = \mathbb{E} \left[D \left(\mathbb{E}[\Theta] \parallel P_{X_0|X_{-k}^{-1}}(\cdot|Y_{-k}^{-1}) \right) \right] + H(\mathbb{E}[\Theta]) - \mathbb{E}[H(\Theta)], \quad (1)$$

where on the right side $\Theta \sim \Pi$ and $Y_{-k}^{-1} \stackrel{i.i.d.}{\sim} \mathbb{E}[\Theta]$.

Corollary III.7 (Relative Entropy with Markovian Law). *For $Q^{\text{LZ},\Pi}$ with $\text{supp}(\Pi) = \mathcal{M}(\mathcal{A})$, and any k^{th} -order Markov law $P \in \mathcal{M}_k$,*

$$\lim_{n \rightarrow \infty} \frac{1}{n} D \left(Q_{X^n}^{\text{LZ},\Pi} \parallel P_{X^n} \right) = \mathbb{E} \left[D \left(\mathbb{E}[\Theta] \parallel P_{X_0|X_{-k}^{-1}}(\cdot|Y_{-k}^{-1}) \right) \right] + H(\mathbb{E}[\Theta]) - \mathbb{E}[H(\Theta)], \quad (2)$$

where on the right side $\Theta \sim \Pi$ and $Y_{-k}^{-1} \stackrel{i.i.d.}{\sim} \mathbb{E}[\Theta]$.

Remark III.8. In the spirit of Remark III.5, the expression in the right sides of (1) and (2) can be equivalently written as

$$D \left(P_Y \parallel P_{X_0|X_{-k}^{-1}} \Big| P_{Y_{-k}^{-1}} \right) + I(\Theta; Y),$$

where: $\Theta \sim \Pi$, $Y \sim \Theta$, $Y_{-k}^{-1} \stackrel{i.i.d.}{\sim} Y$, $P_{X_0|X_{-k}^{-1}}$ is that associated with the $P \in \mathcal{M}_k$, and here we assume the notation of [11] for relative entropy between conditional distributions (kernels) averaged with respect to another distribution.

[SDHW26]:

distributional properties

Corollary IV.11. *Applying Theorem IV.10 with single-sequence event A such that $A_1 = \dots = A_r = \mathcal{M}(\mathcal{A})$, we obtain*

$$\frac{1}{n} \sum_{i=1}^n \delta_{X_i^{i+r-1}}(a^r) \xrightarrow{a.s.} \prod_{t \in [r]} \mathbb{E}[\Theta[a_t]] \quad \text{as } n \rightarrow \infty, \quad \forall a^r \in \mathcal{A}^r.$$

In words, the r th-order empirical distribution of the realized sequence from the LZ78 source converges, almost surely, to the law of r i.i.d. drawings from the PMF $\mathbb{E}[\Theta]$.

intuition

- the true model keeps track of the tree and the location of the process in it
- most of the time the process is in nodes visited a growing number of times so true model quickly learns their realized theta
- then “merely” need show the empirical distribution of the thetas is what one would expect
- any finite memory machine loses sync of where the process is in the tree and since things are locally iid its log loss is as good as that for predicting an iid source

[SDHW26]:

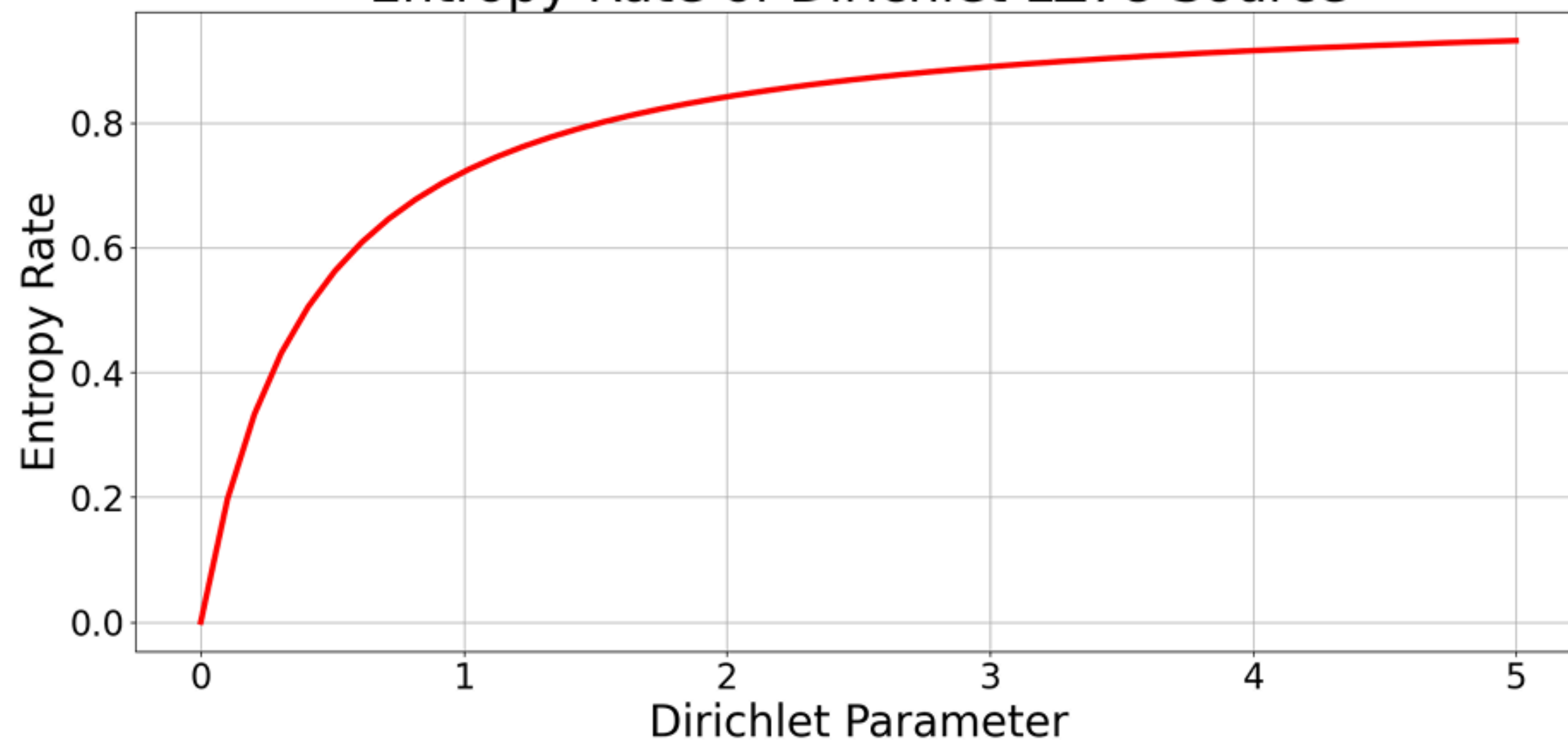
question:

what are the distributional and entropic properties of a process generated from the LZ78 source?

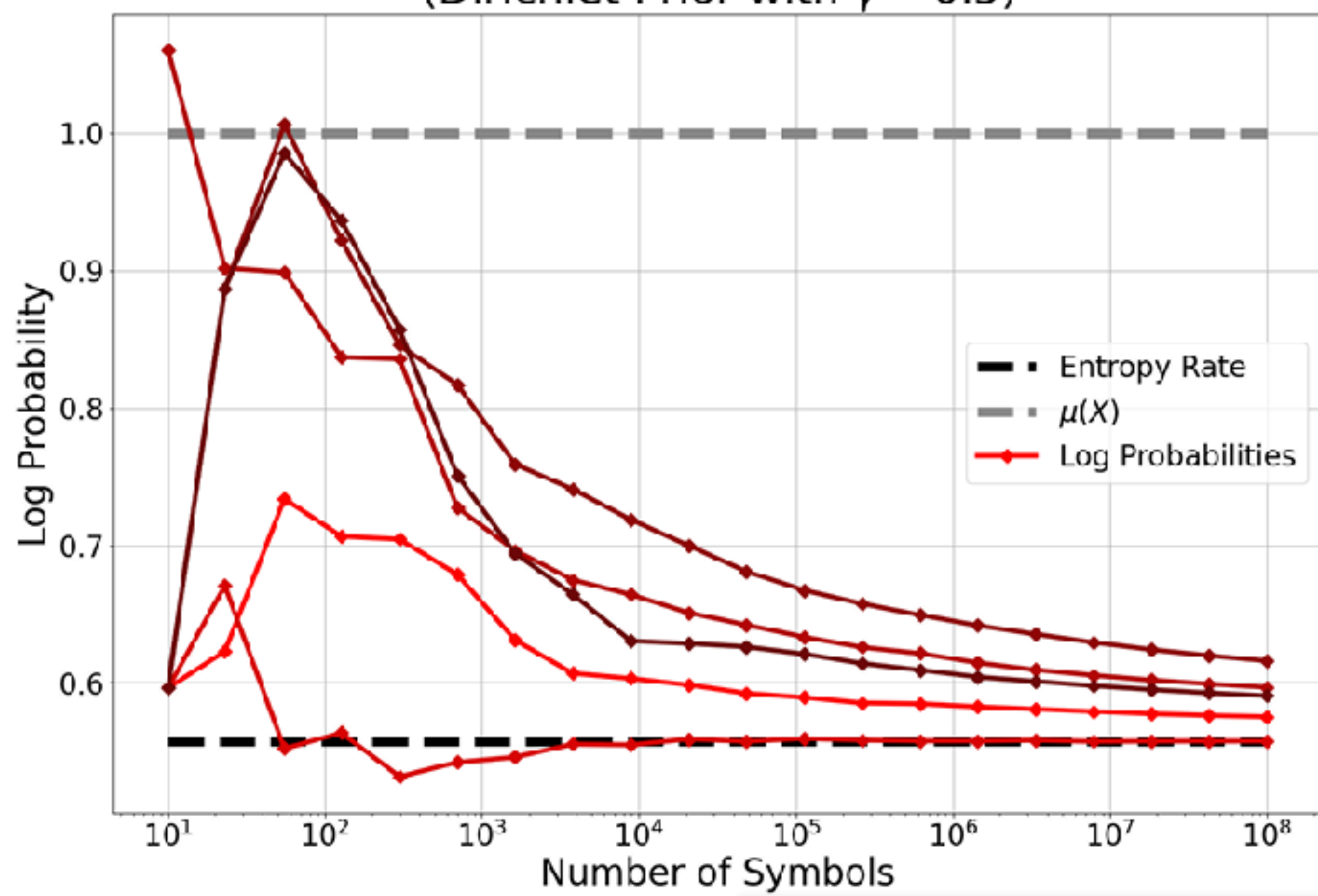
answer: the “LZ source”

- is non-stationary
- but ergodic
- is asymptotically “locally iid”
- is globally far from it
- has entropy rate $<$ finite-state compressibility
- has the Shannon-McMillan-Breiman property
- divergent relative entropies

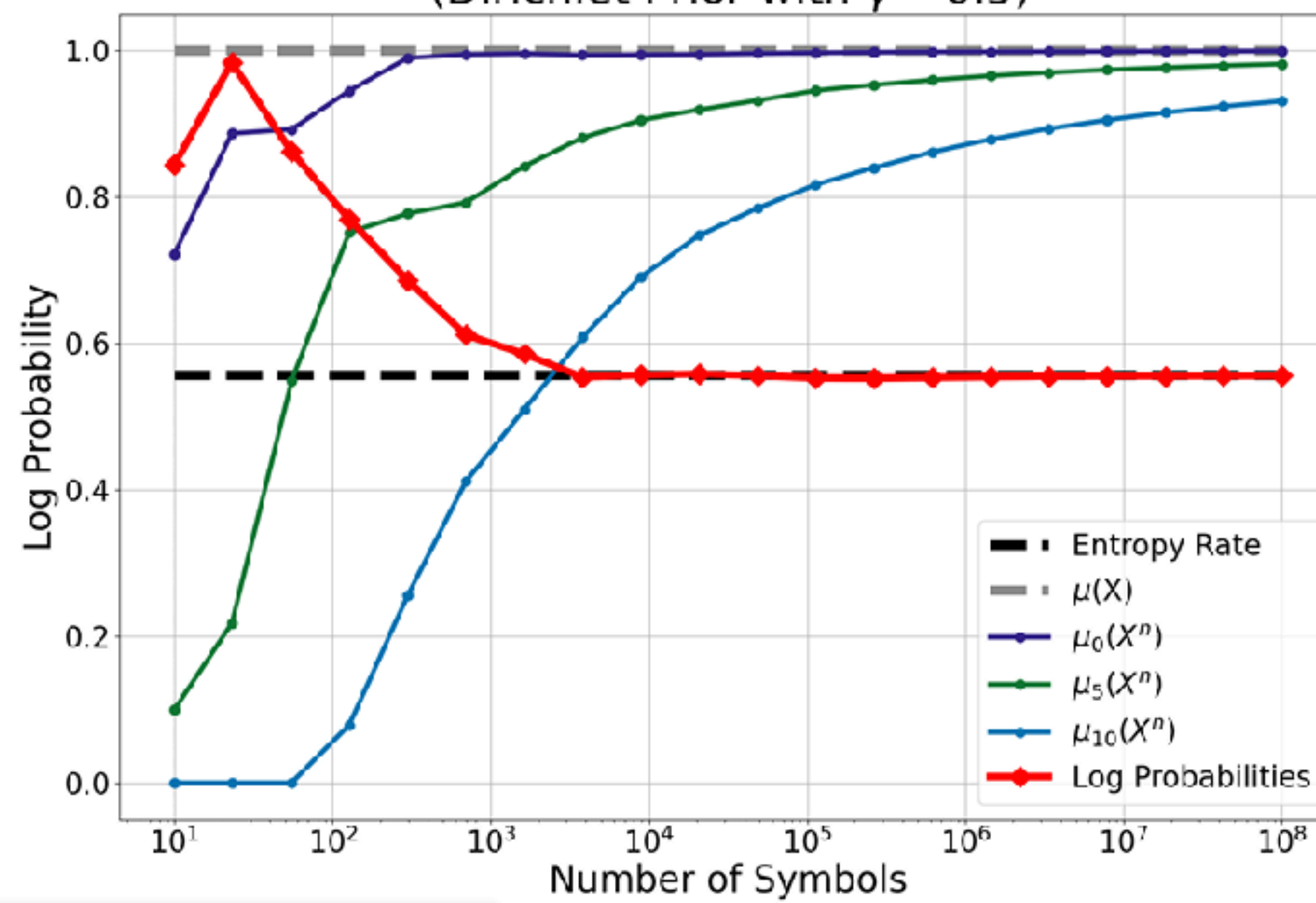
Entropy Rate of Dirichlet LZ78 Source



Log Probability of Sequences from LZ78 Source (Dirichlet Prior with $\gamma = 0.5$)



Log Probability of Sequences from LZ78 Source (Dirichlet Prior with $\gamma = 0.5$)



why is this interesting?

- provides a highly non-trivial ergodic information source with very different local vs global behavior

cf. upcoming international symposium on information theory (ISIT 2026):

- can be useful for testing deep nets
- might be a ‘universal data source for training’

The LZ78 Source and its Application to Evaluating In-Context Learning

Naomi Sagan
Electrical Engineering
Stanford University
Stanford, CA, USA
nsagan@stanford.edu

Amir Dembo
Statistics, Mathematics
Stanford University
Stanford, CA, USA
adembo@stanford.edu

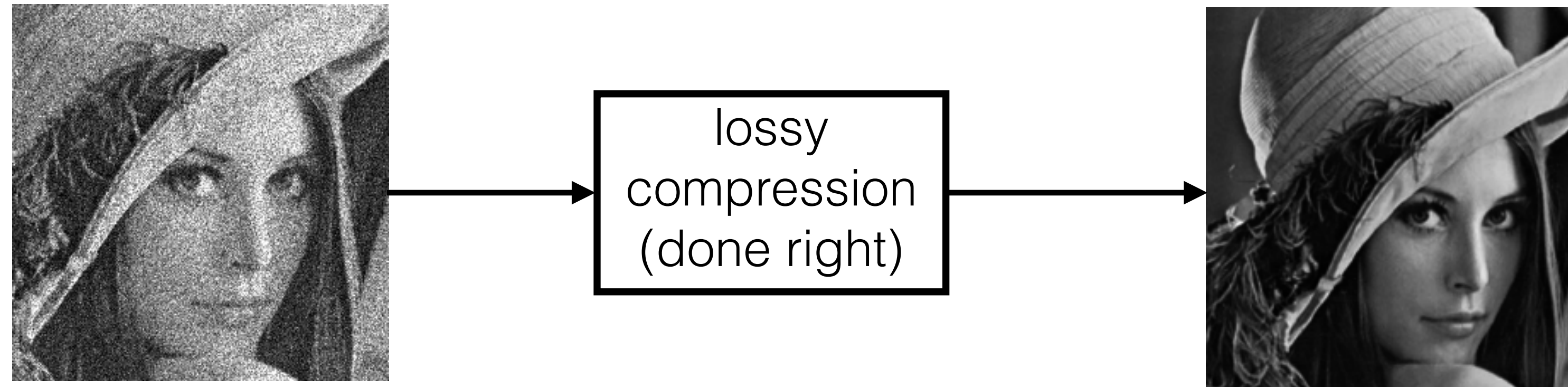
Matthew Ho
Electrical Engineering
Stanford University
Stanford, CA, USA
mattho@stanford.edu

Tsachy Weissman
Electrical Engineering
Stanford University
Stanford, CA, USA
tsachy@stanford.edu

- can be harnessed for lossy compression

speaking of lossy compression

denoising via lossy compression



- “Occam filters”
- “Compresstimation”
- “Kolmogorov sampler”
- etc.

renewed recent interest due to prevalence of good “universal” lossy compressors

lossy compression of noisy signals for denoising (cont.)

Effect of lossy compression of quality scores on variant calling

[Idoia Ochoa](#) ✉, [Mikel Hernaez](#), [Rachel Goldfeder](#), [Tsachy Weissman](#), [Euan Ashley](#)

Briefings in Bioinformatics, Volume 18, Issue 2, March 2017, Pages 183–194,

Compression Detects Changes in Spiking Neural Data from Cortical Lesions

Alice Tor¹, Yuxin Wu¹, Stephen E Clarke^{2,3}, Lisa Yamada^{2,1,3}, Tsachy Weissman¹, Paul Nuyujukian^{2,3,1,4}, for the Brain Interfacing Laboratory

Lossy Compression of Noisy Data for Private and Data-Efficient Learning

Berivan Isik and Tsachy Weissman, *Fellow, IEEE*

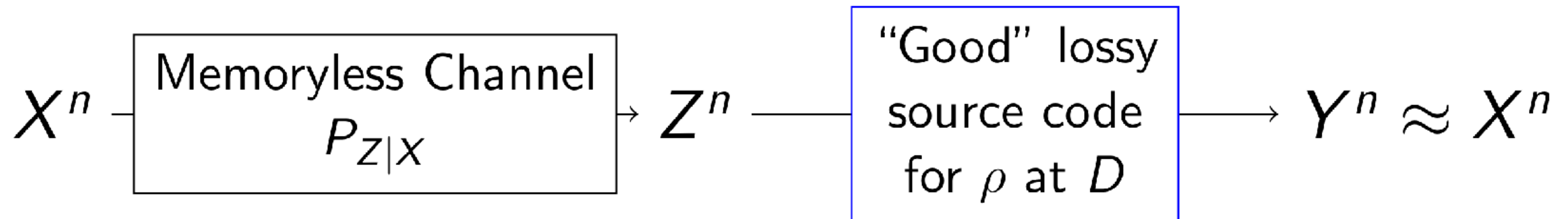
less storage, faster processing/learning, boosted accuracy in the downstream, enhanced privacy when it matters, ...

[SOW26]: “The Performance of Compression-Based Denoisers”,

Dan Song, Ayfer Özgür and T. Weissman,

to appear in Information Theory Transactions

<https://arxiv.org/abs/2512.14539>



- ▶ Source X^n from fixed stationary ergodic source X .
- ▶ Channel $P_{Z|X}$ are fixed and known.

“Good” Codes

Definition

A sequence of codes (\mathcal{C}_n, ϕ_n) is *good* at distortion D if asymptotically achieves a point on the rate distortion curve:

$$\limsup_{n \rightarrow \infty} \mathbb{E}[\rho_n(Z^n, \phi_n(Z^n))] \leq D.$$

$$\limsup_{n \rightarrow \infty} \frac{1}{n} \log |\mathcal{C}_n| \leq \lim_{k \rightarrow \infty} R(Z^k, D)$$

We will take $Y^n = \phi_n(Z^n)$.

[SOW26]:

Empirical Distribution

Let:

$$J \sim \text{Unif}([n - 2k - 1]), J \perp (X^n, Z^n, Y^n).$$

Define:

$$X^n = (X_1, \dots, X_{J-1}, X_J, \dots, X_{J+2k}, X_{J+2k+1}, \dots, X_n)$$

$$Z^n = (Z_1, \dots, Z_{J-1}, Z_J, \dots, Z_{J+2k}, Z_{J+2k+1}, \dots, Z_n)$$

$$Y^n = (Y_1, \dots, Y_{J-1}, \underbrace{Y_J, \dots, Y_{J+2k}}_{(\tilde{X}^{(n,k)}, \tilde{Z}^{(n,k)}, \tilde{Y}^{(n,k)})}, Y_{J+2k+1}, \dots, Y_n).$$

Let $Q^{(n)}$ be the probability measure such that

$$(\tilde{X}^{(n,k)}, \tilde{Z}^{(n,k)}, \tilde{Y}^{(n,k)}) \sim Q_{X_{-k}^k, Z_{-k}^k, Y_{-k}^k}^{(n)}.$$

[SOW26]:

Theorem 2.2 *Suppose Y^n is the output of a good compressor for Z^n for distortion*

$$\rho(z, y) = -\log P_{Z|X}(z|y)$$

at distortion level

$$D = H(Z|X).$$

Then, for every k ,

$$Q_{Z^k, Y^k}^{(n)} \xrightarrow{n \rightarrow \infty} P_{Z^k, X^k}$$

in words:

'good' lossy compression of the noisy signal gives a 'sample from the posterior'

[SOW26]:

Characterization of Loss

Theorem (Informal)

For any choice of bounded loss, the loss between X^n and Y^n achieved by a sequence of good codes for

$$X^n \text{ --- } Z^n \text{ --- } Y^n$$

approaches the loss between X_0 and X'_0 , where

$$P_{X'_0|Z} = P_{X_0|Z},$$

$$X_0 \text{ --- } Z \text{ --- } X'_0.$$

in words: an *independent* sample from the posterior

Exact Characterization of Loss

Theorem

Suppose (\mathbf{X}, \mathbf{Z}) are double-sided mixing, $P_{Z|X}$ has invertible channel matrix, and $\mathcal{X}, \mathcal{Z}, \mathcal{Y}$ are finite.

Choosing distortion $\rho(z, y) = -\log p_{Z|X}(z | y)$ and distortion level $H(Z | X)$, for any bounded loss $\Lambda : \mathcal{X} \times \mathcal{Y} \rightarrow [0, \Lambda_{\max}]$,

$$\lim_{n \rightarrow \infty} \mathbb{E}[\Lambda_n(X^n, Y^n(Z^n))] = \mathbb{E}_{Z_{-\infty}^{\infty}} \left[\mathbb{E}_{(X', Y') | Z_{-\infty}^{\infty}} [\Lambda(X', Y')] \right], \quad (14)$$

where

$$(X', Y') \sim \left(P_{X_0 | Z_{-\infty}^{\infty}} \right)^2. \quad (15)$$

- ▶ Loss is **exactly characterized** (the limit exists).
- ▶ **Independent** samples from the posterior $P_{X_0 | Z_{-\infty}^{\infty}}$.

(note scheme is oblivious to loss function)

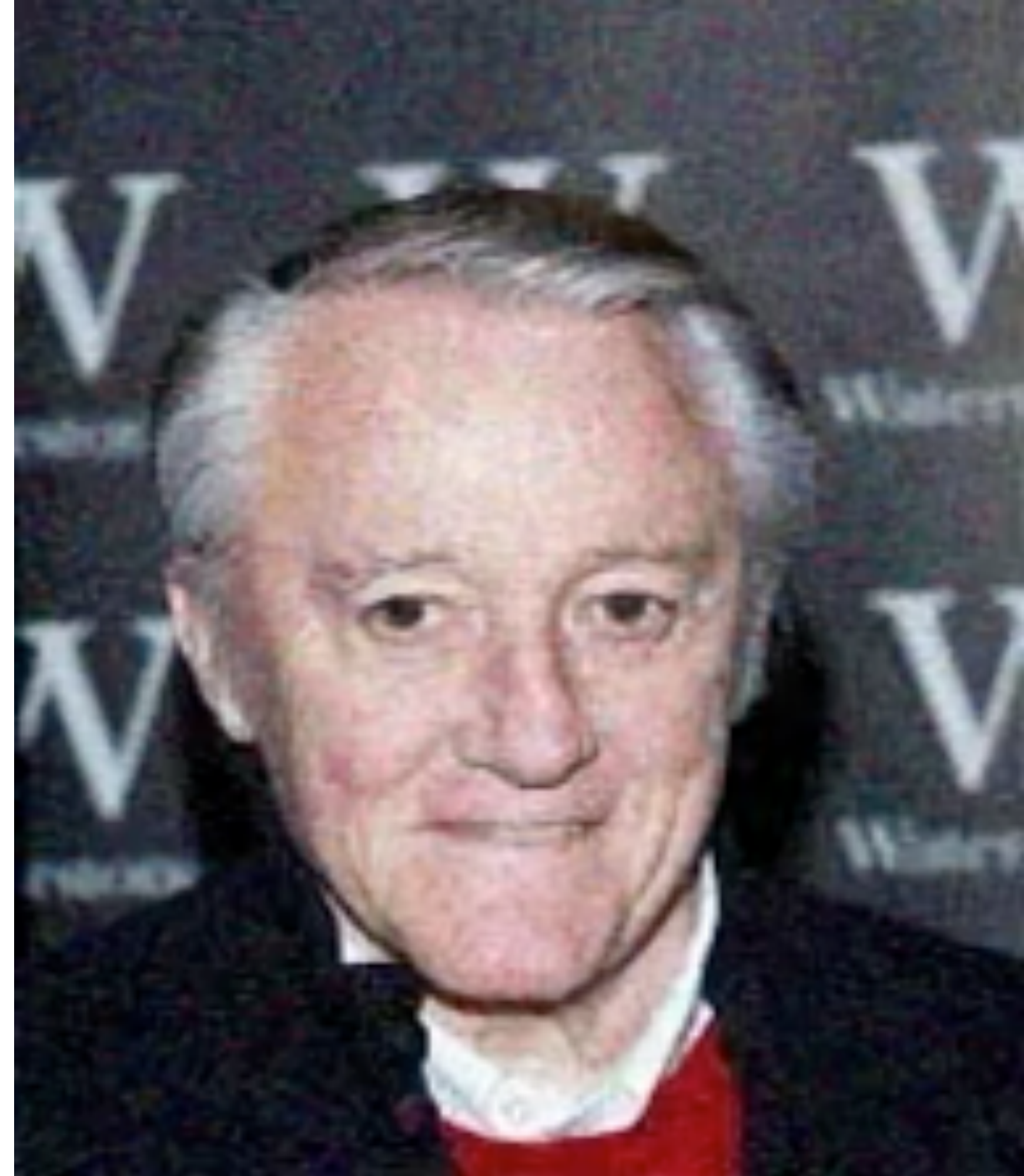
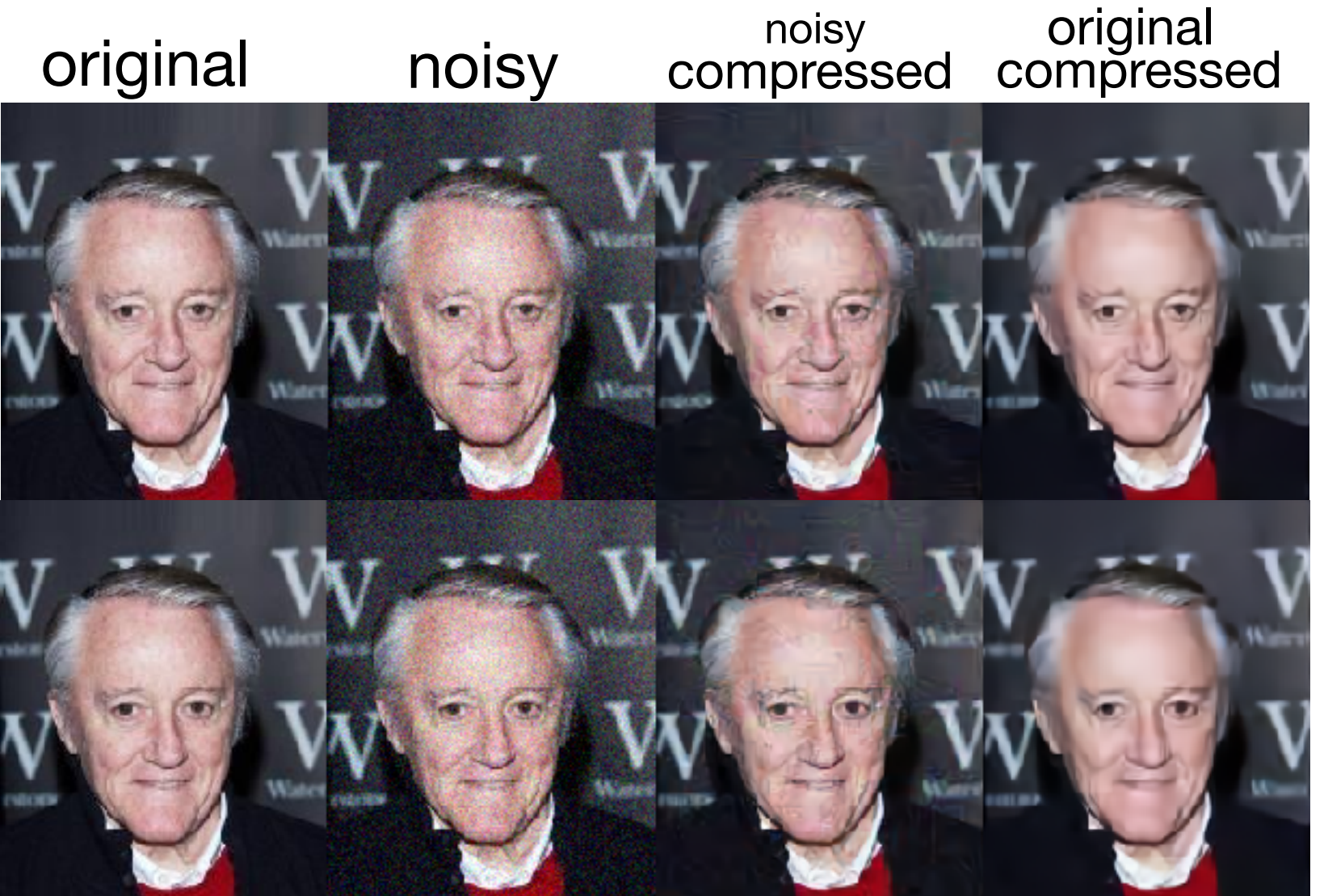
proof elements

- the r-d achieving test channel of the noisy source is the true posterior
- empirical distribution of good codes \sim the r-d achieving test channel
- new Markov lemma \Rightarrow a conditionally independent sample from the posterior

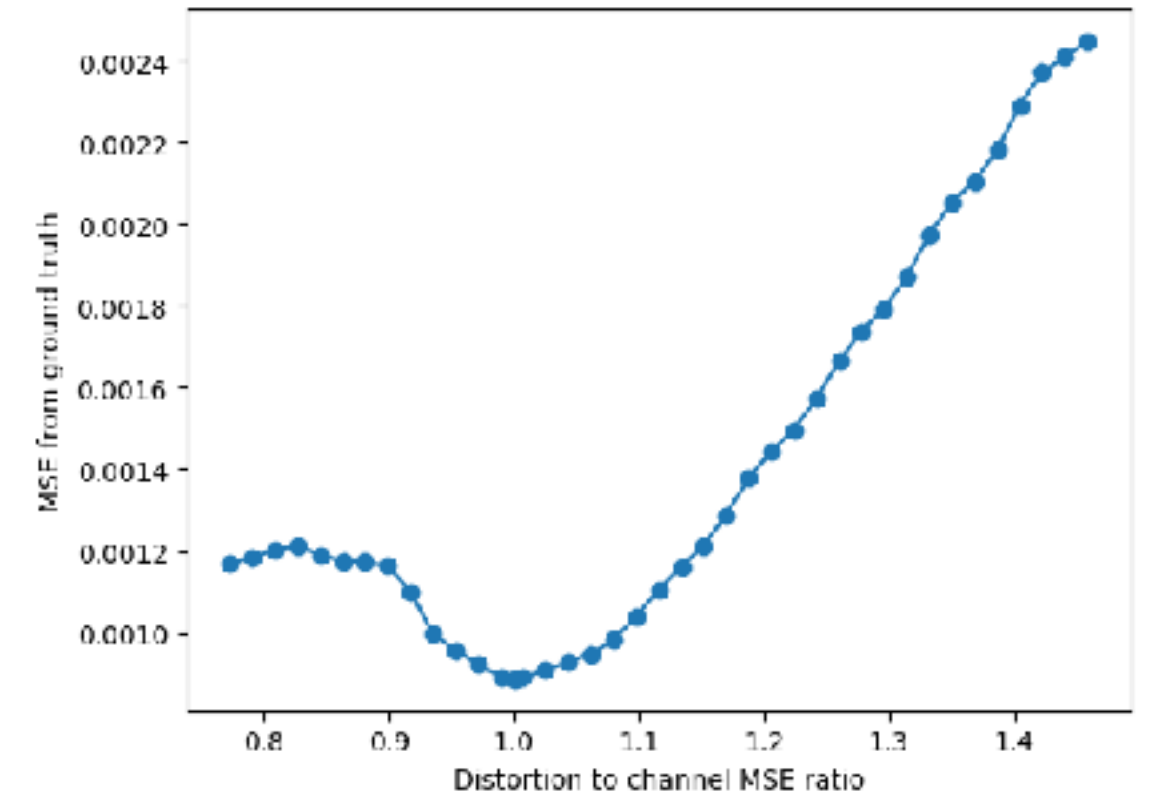
benefits

- generalizes, strengthens, and bridges gaps in the literature
- oblivious to the loss function under which denoising is judged
- yields schemes with “perfect perception” a la “rate-distortion-perception”

AWGN
+
squared
error
distortion



(also yields
state of the art
rate-distortion-
perception
performance)



Unsupervised Audio Denoising via Neural Lossy Compression

- **Motivation:** Neural codecs for audio are increasingly prevalent which much success. Neural compressors offer better rate-distortion performance than traditional compressors. Since audio is 1-dimensional, we can obtain a good lossy compressor with a small-scale network.
- **Approach:** We train a scale-hyperprior neural compressor on **noisy audio** samples to construct a good lossy code for the given noise mechanism. The learned compressor can be applied to the target noisy signal, which upon reconstruction, gives us a denoised signal. Thus, we obtain **joint audio compression and enhancement**.
- **Unsupervised Benefit:** There have been recent attempts at this joint problem, such as SoundStream [1] and ADNAC [2]. However, these approaches require access to clean-noisy audio pairs, which may not always be viable. Our approach only requires access to the noisy data.
- In addition to compression and denoising, we have a positive byproduct of bandwidth reduction.

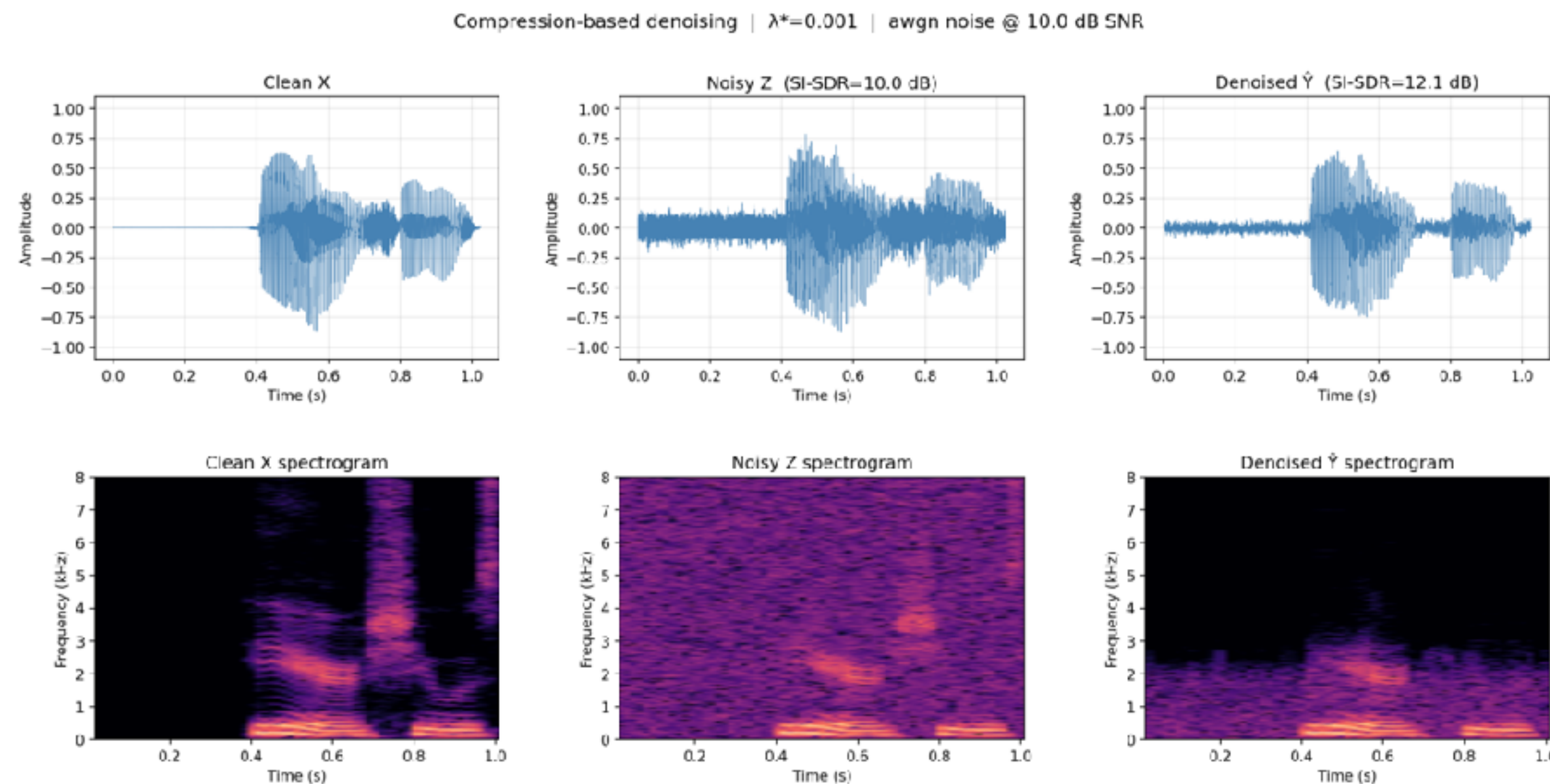


Figure: Clean, noisy, and denoised signals at 10dB SNR AWGN noise.

- [1] "SoundStream: An End-to-End Neural Audio Codec," N. Zeghidour et al.
 [2] "ADNAC: Audio Denoiser using Neural Audio Codec," D. Jimon et al.

Results and Takeaways:

- **Compression and Denoising Performance:**
 - At 10dB SNR, we get ~4dB gain and ~12x compression.
 - At 5dB SNR, we get ~6dB gain and ~20x compression.
- **Classical Baselines:** We compare the denoising performance of our compressor to classical methods such as lowpass filtering, spectral subtraction, and Wiener filtering. In terms of SI-SDR, our approach outperforms these methods
- **Neural Audio Codecs:** We look at two off-the-shelf neural audio codecs used in practice – EnCodec and DAC – and compare their rate-denoising performance to our approach
 - EnCodec and DAC are neural codecs trained on, and meant for, clean audio data
 - Our approach gets a better denoising gain than these codec, while achieving a comparable bitrate regime
- **Specialised Audio Denoisers:** We also compare our denoising performance to Facebook Demucs [1], which is a specialized audio denoiser trained in a supervised way, and outperforms our denoising capabilities. However, it offers no compression, and should be treated as a denoising benchmark

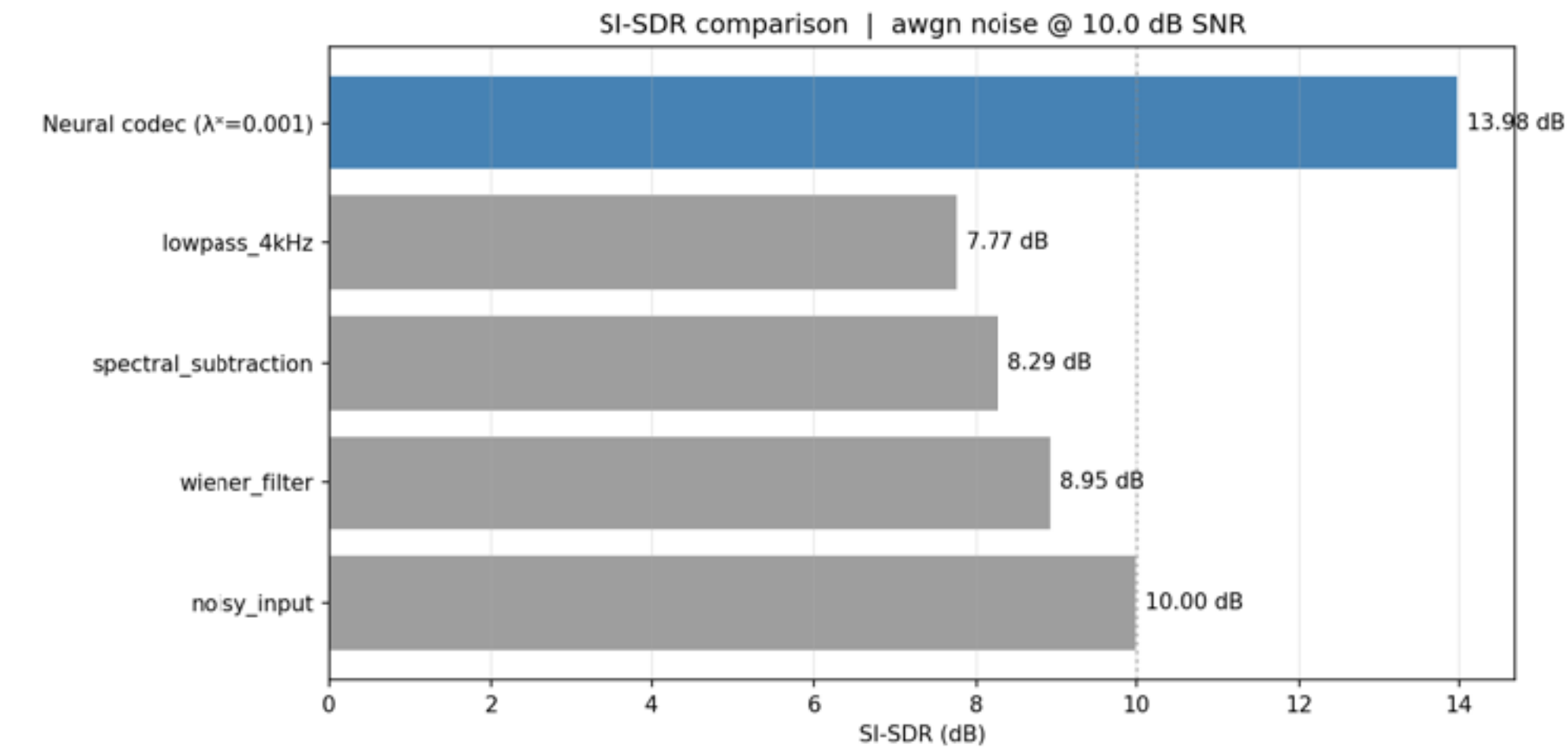


Figure: Denoising performance at 10dB SNR against classical baselines.

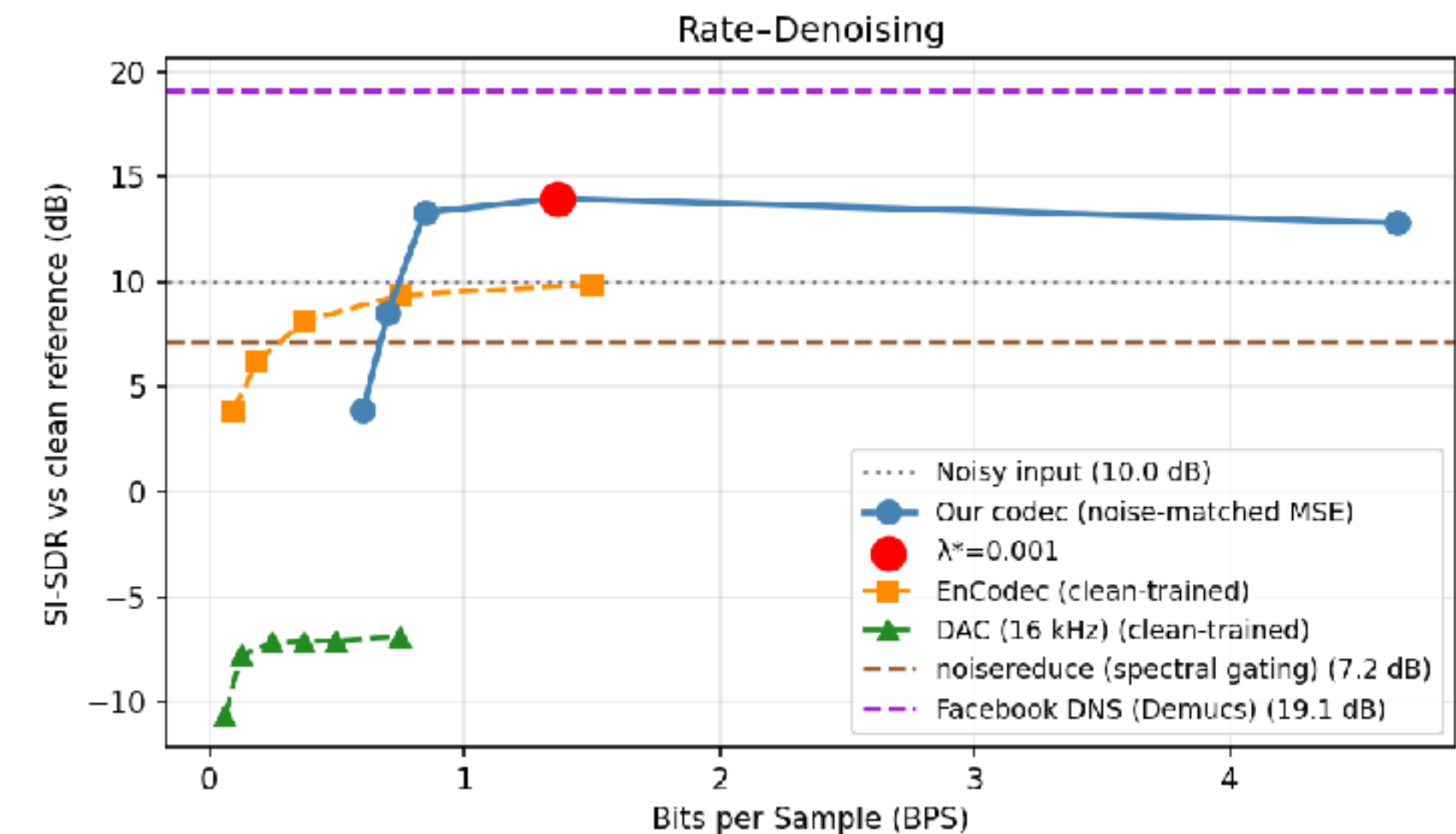


Figure: Comparing rate-denoising performance at 10dB SNR against off-the shelf methods.

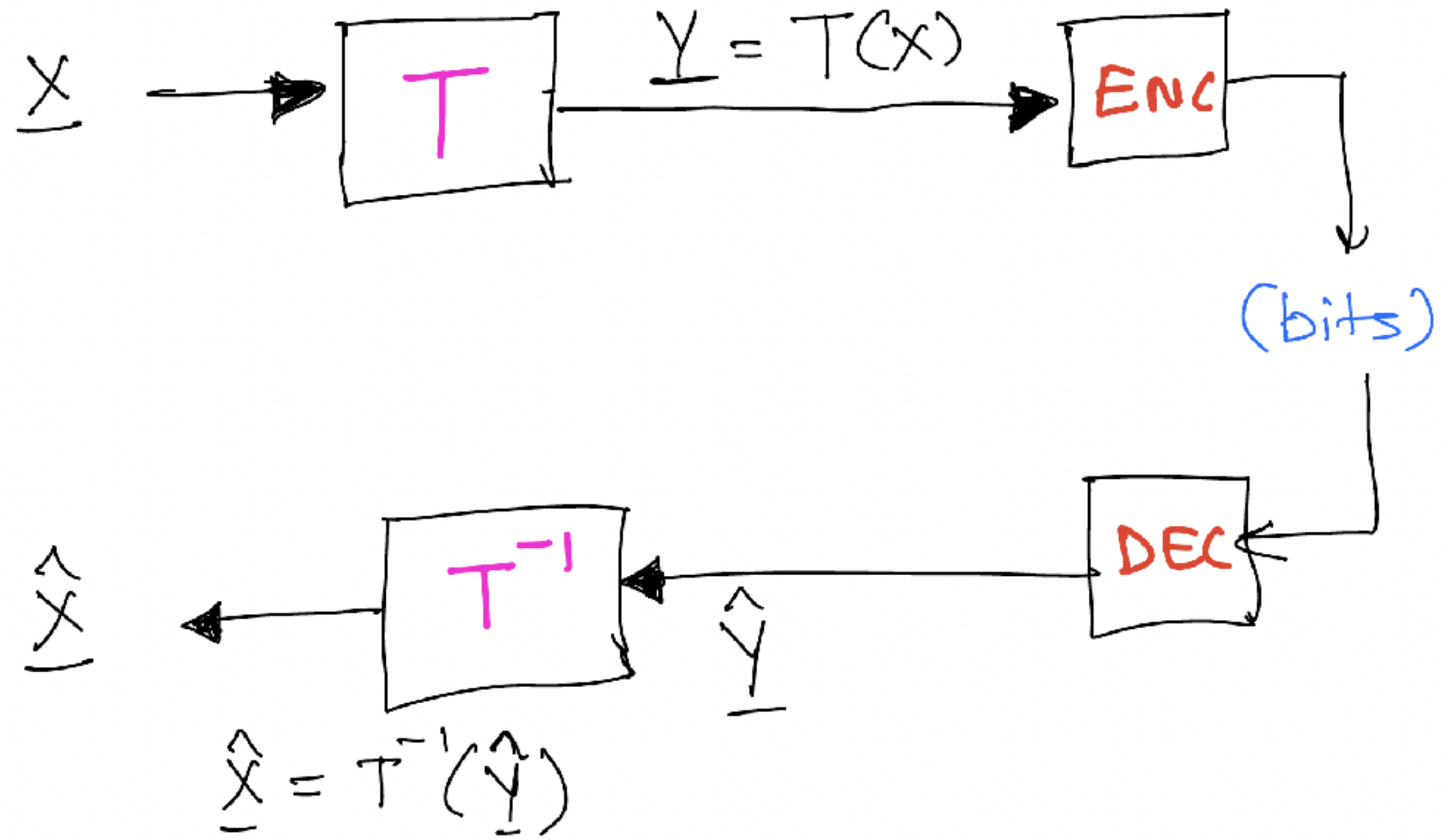
[1] “Real Time Speech Enhancement in the Waveform Domain,”
A. Défossez et al.

C. Ding, A. Rao Gorle, J. Jeong, N. Sagan and T. Weissman,
"Information-computation trade-offs in non-linear transforms",
Philosophical Transactions of the Royal Society A, 2026
[Ding et al. 26]

among other highly non-linear transforms, considers:

- neural
- textual
- LZ

transform coding 101



elements of a good transform

- basis functions have biological/physical/conceptual significance
- sparsity
- transform coefficients related simply (e.g. independent, uncorrelated, etc.)
- invertibility (information lossless)
- smoothness of forward and inverse transform (w.r.t. relevant metrics)
- manageable complexity (both forward and inverse transform)
- induced benefit: yields a wieldy and natural model for your data

quest for a new transform

- T. Weissman, “Toward textual transform coding”, IEEE BITS, Vol 3, Issue 2, June 2023, [W23]
 - what resonates with us?
 - what is a ‘thing’?
 - what lights up our (biological) neural nets?

what are the right “harmonics”?

words

- are what we've invented/evolved to refer to and describe concepts/things that resonate with/are meaningful to us

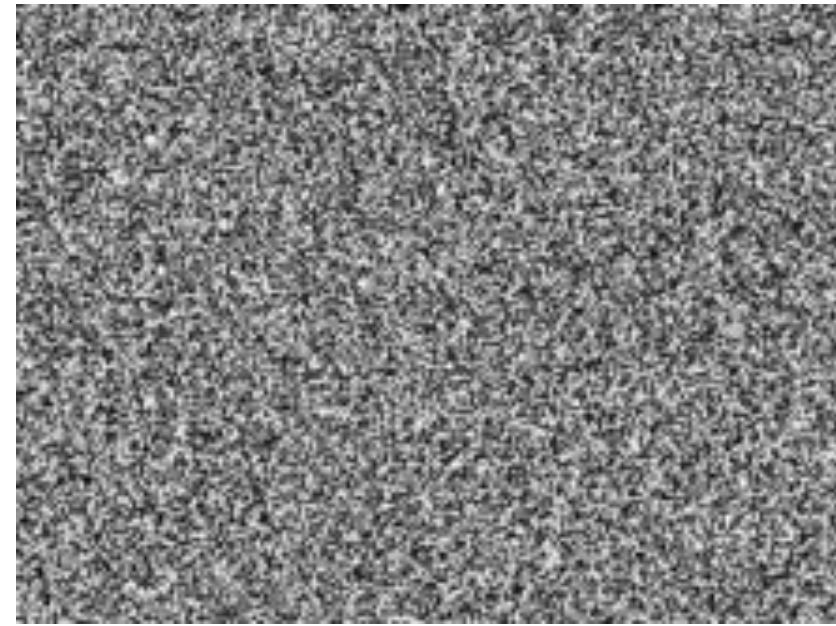
analogy toward a “textual transform”

Fourier transform	Textual transform
harmonics	objects for which we have words
frequencies	words
amplitude/phase of harmonics	size/location/orientation of objects
precision of amplitude/phase	precision of size/location
frequency resolution	number of words

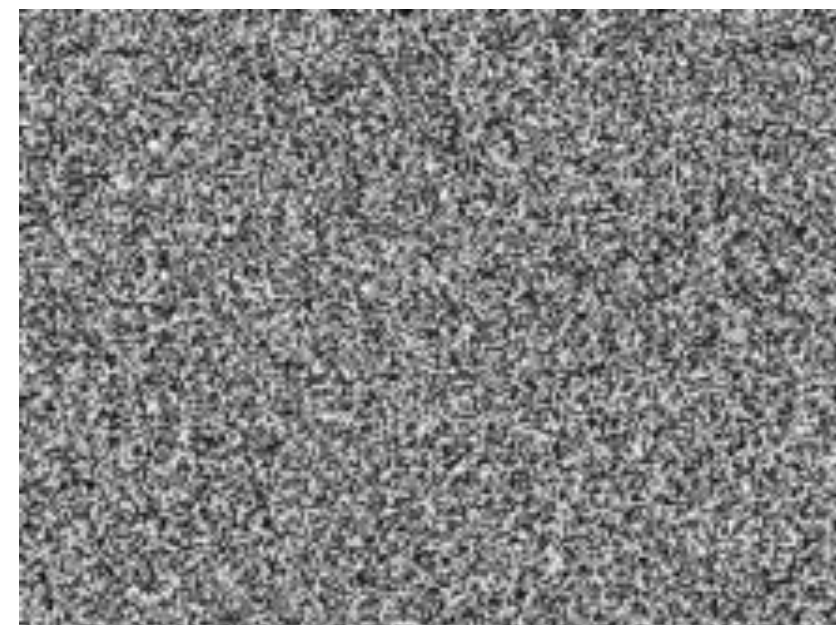
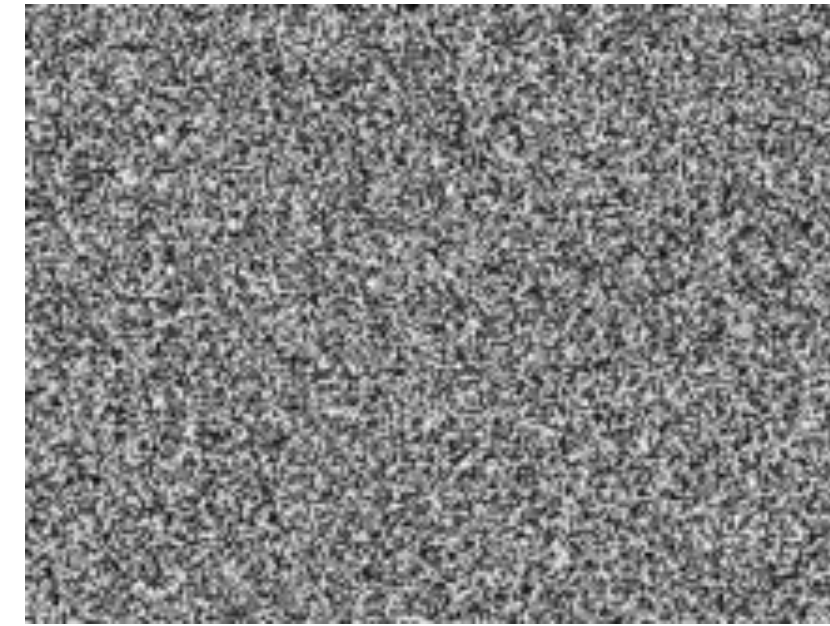
(some cartoonish sketches of)

examples for what such a transform might look like

white noise



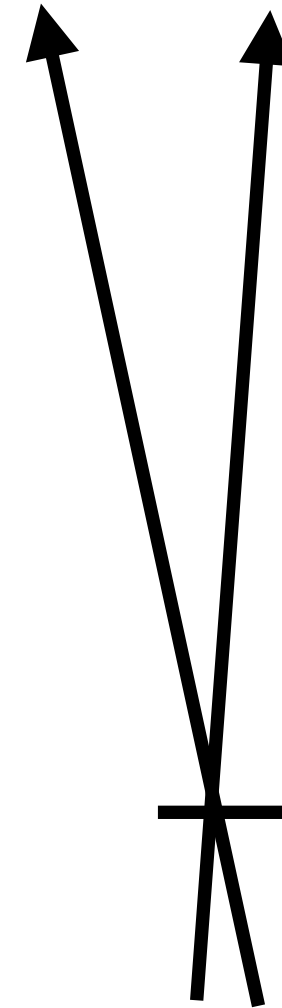
FFT, etc.



Textual Transform



size,
location,
orientation, etc.



white noise

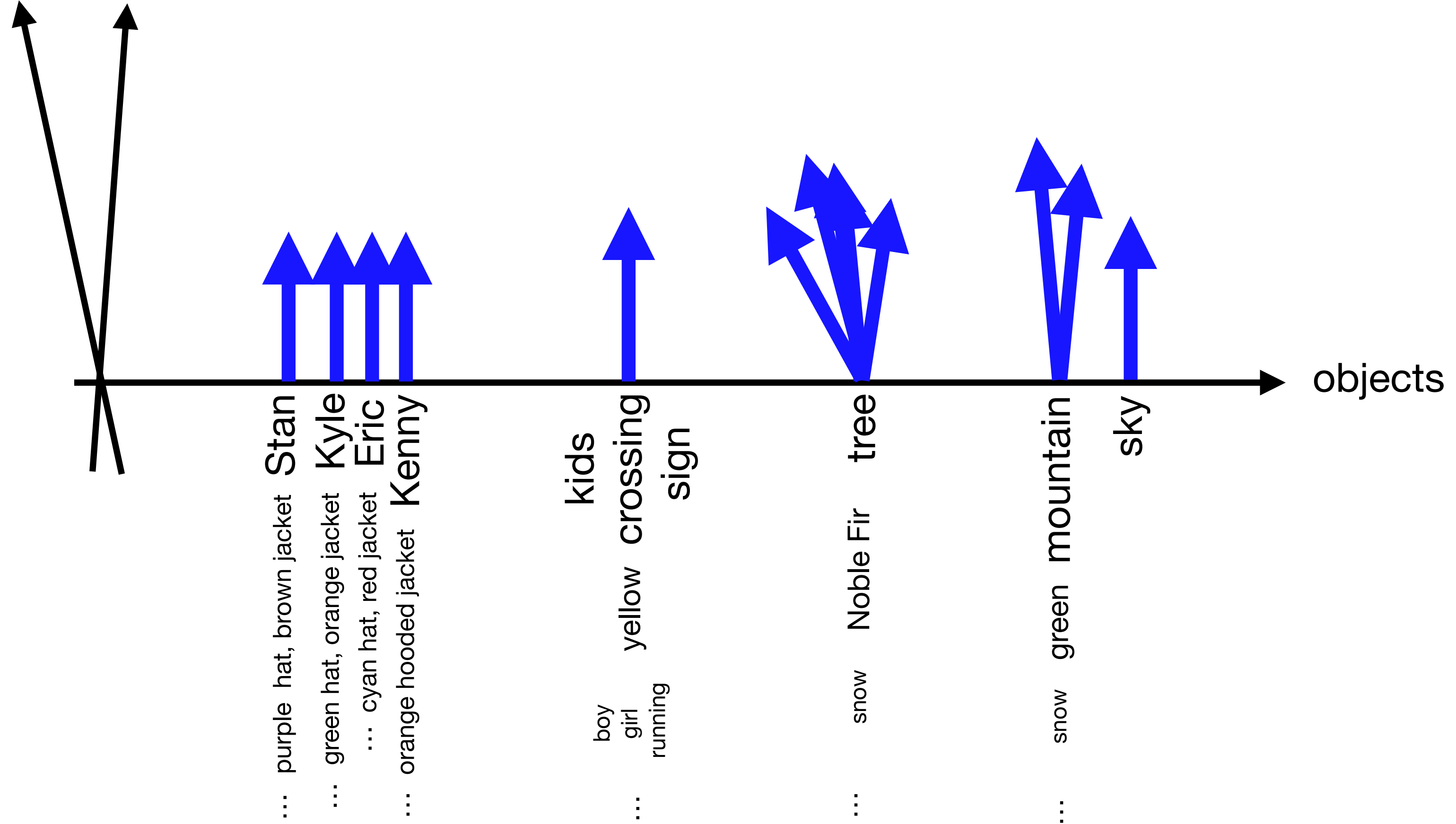


objects

South Park image



size,
location,
orientation, etc.




representation of an image in the textual transform domain

- how many objects (up to prescribed number L for cognitive load)
- list of what they are (one word each)
- their sizes, locations, orientations (at prescribed resolution R)
- more words of description for each and how they relate to each other (at prescribed resolution W)

in other words:

**essentially what a verbal description would look like
(confined to a given budget of words/sentences)**

textual transformers

- transformers (Large Language and visual models), diffusion models, etc.
- computationally effective textual transforms
- humans 
- ongoing: text emerging from extreme lossy compression

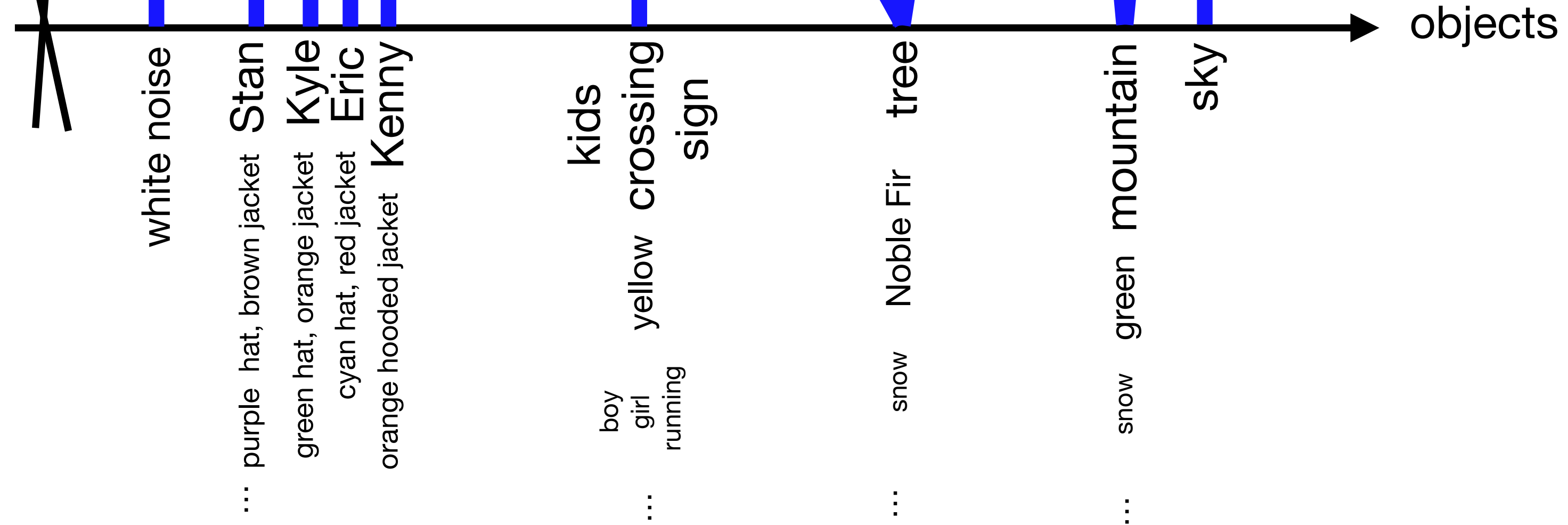
a bit on denoising and lossy compression via the textual transform:

noisy

South Park image



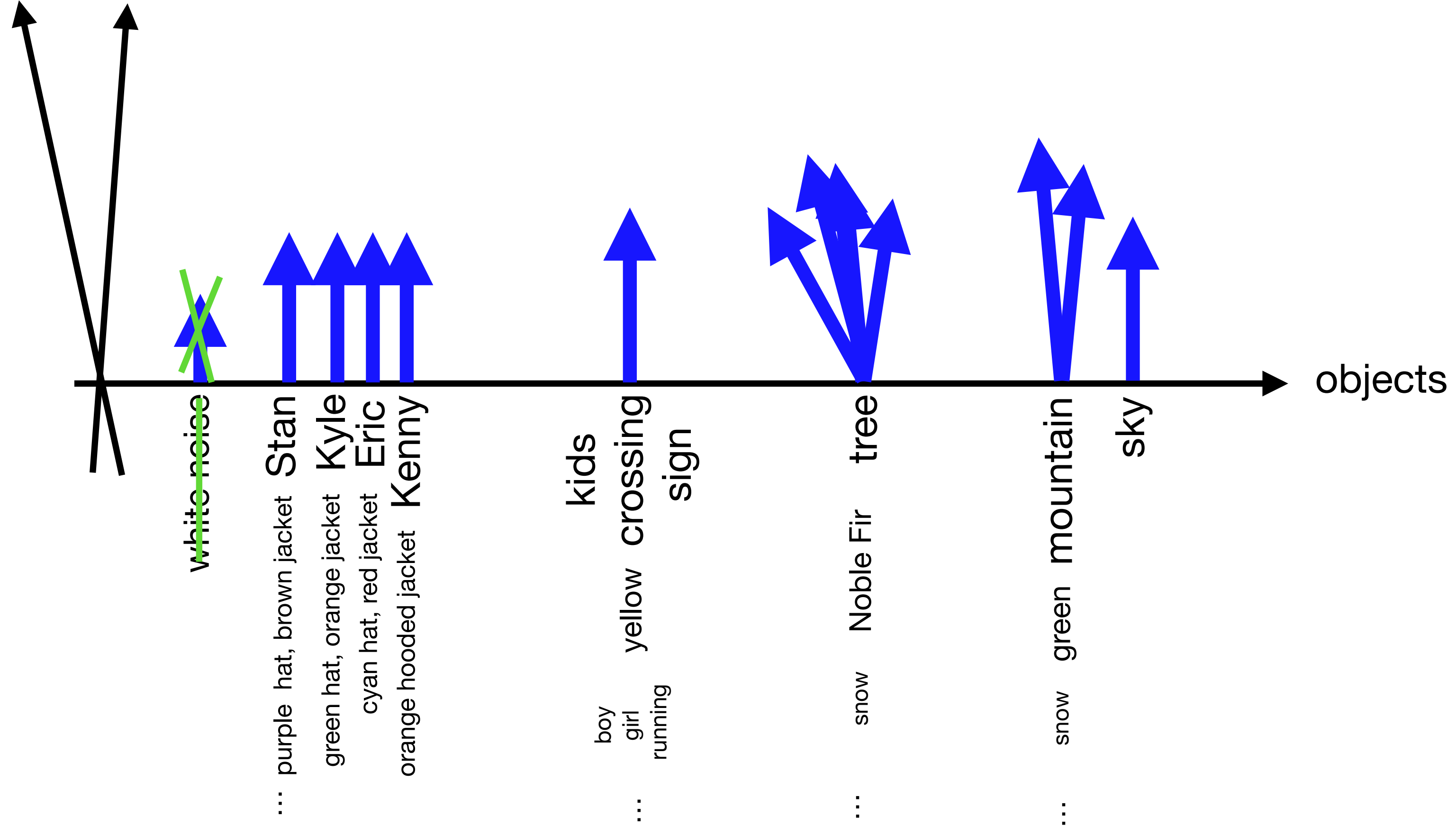
size,
location,
orientation, etc.



denoising the South Park image



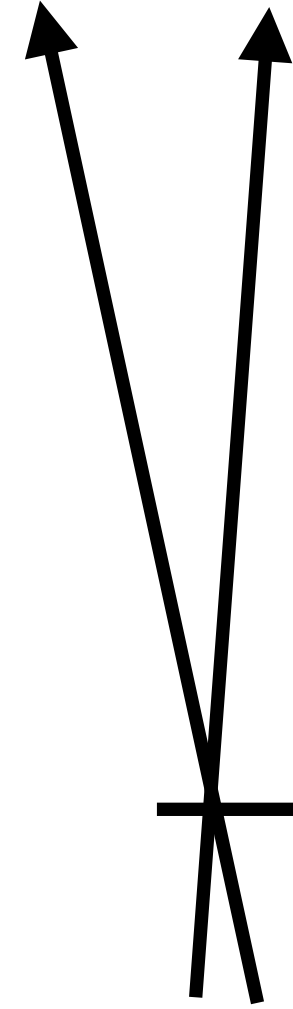
size,
location,
orientation, etc.



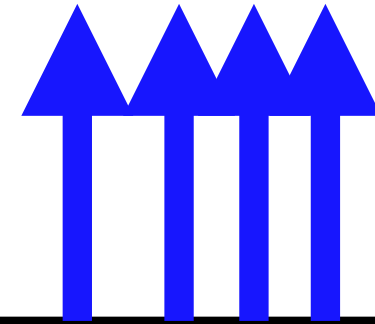
denoised
South Park image



size,
location,
orientation, etc.



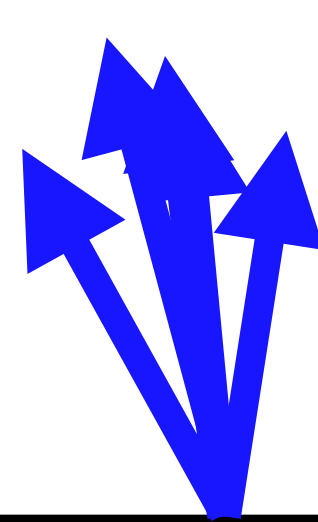
... purple hat, brown jacket
Stan
... green hat, orange jacket
Kyle
... cyan hat, red jacket
Eric
... orange hooded jacket
Kenny



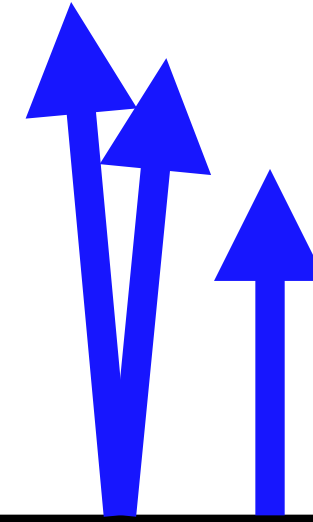
...
kids
... boy
girl
running
yellow crossing
sign



...
snow
Noble Fir
tree



...
snow
green
mountain
sky



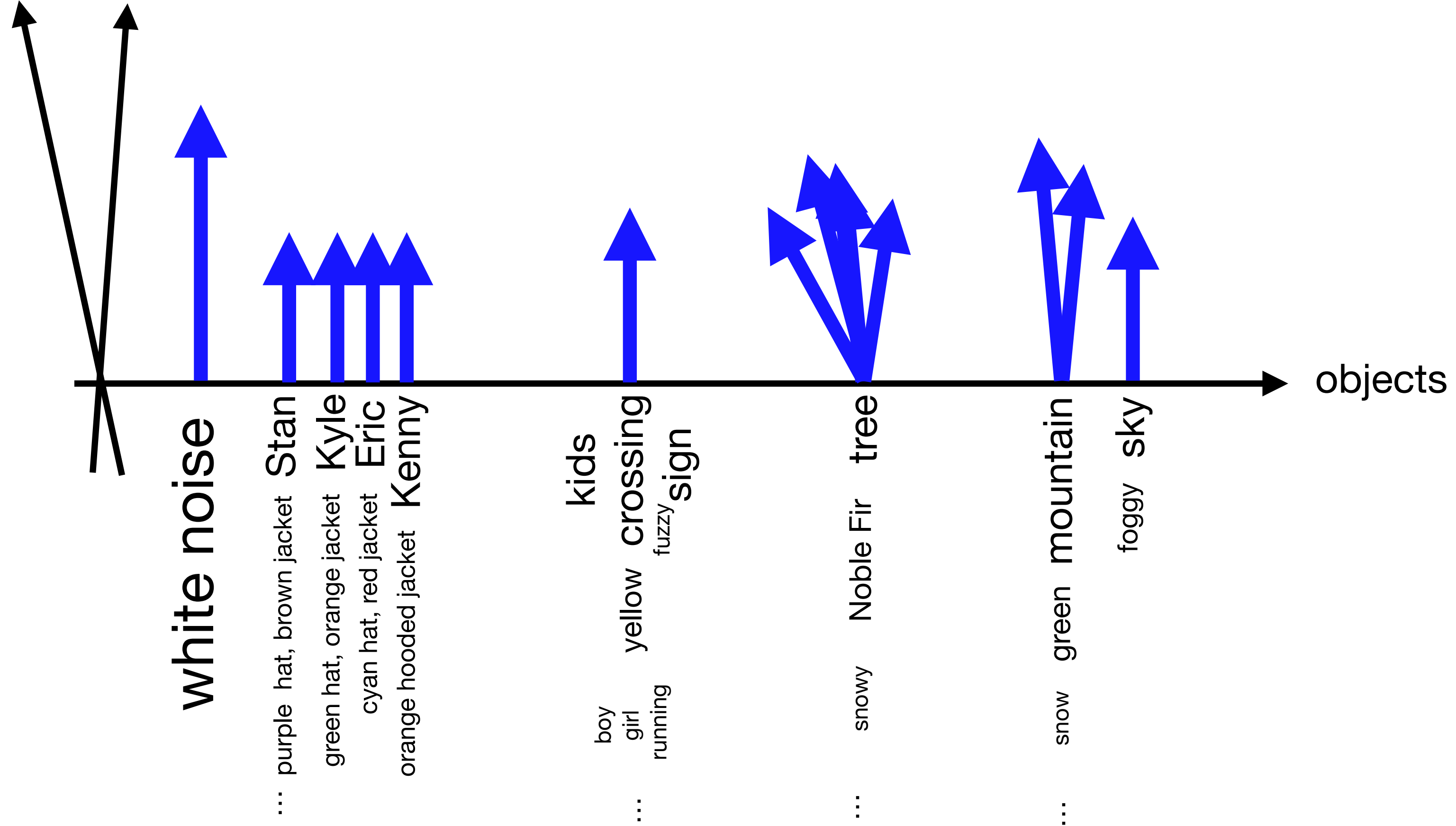
objects

noisier

South Park image



size,
location,
orientation, etc.

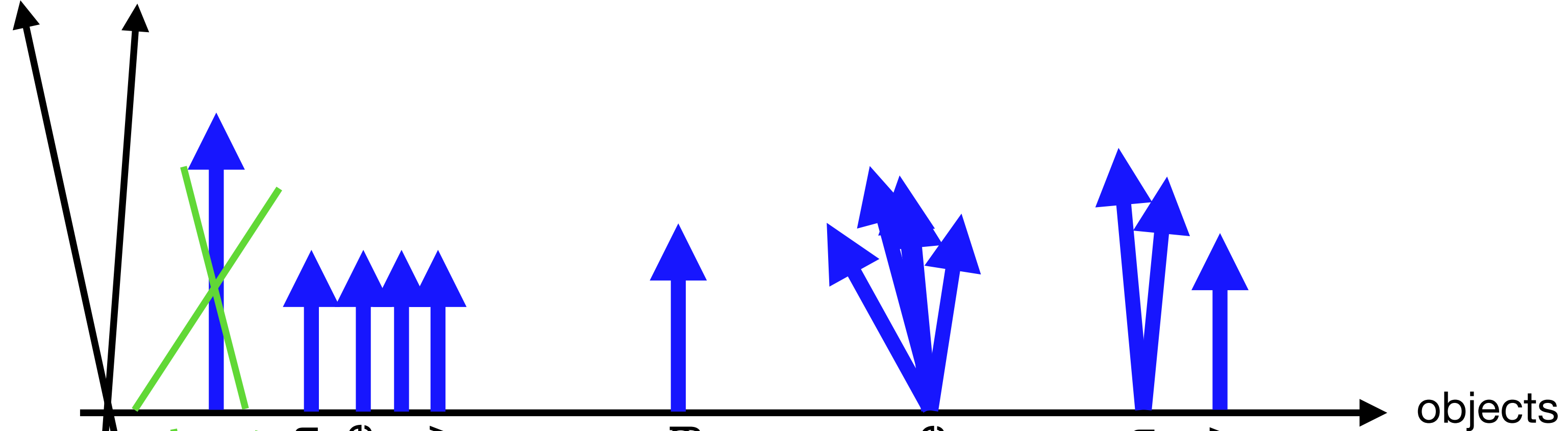


denoising the noisier

South Park image



size,
location,
orientation, etc.



~~white noise~~

... purple hat, brown jacket Stan
green hat, orange jacket Kyle
cyan hat, red jacket Eric
orange hooded jacket Kenny

... kids
boy yellow crossing
girl running ~~fuzzy~~ sign

... ~~snowy~~ Noble Fir tree

... ~~snow~~ green mountain
~~foggy~~ sky

(discrete denoising
in the textual domain)

[Ding et al. 26]: **goes way beyond the cartoons to yield state of the art lossy
compression and denoising**

via the textual transform + visual sketches

some takeaways

- NNs, transformer and diffusion based architectures have shown us what's achievable
- existing compressors such as LZ hint that NNs are highly suboptimal from a compute perspective
- “plenty of room at the bottom” for new architectures
- at compute regimes \gg compressors and \ll NNs
- need go back to the drawing board
- seek more inspiration from humans and other (not necessarily neuronal) biological systems
- new architectures could give rise to new types of intelligence

thanks