

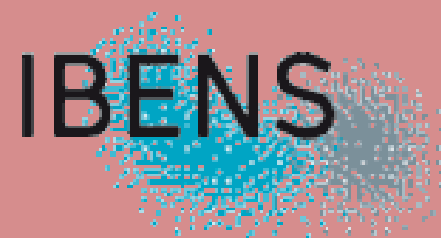
Scaling self-supervised representation learning for symbolic piano performance

Louis Bradshaw, Honglu Fan, ... Stella Biderman, Simon Colton

Presented at



06.01.2026 Clemence Reda



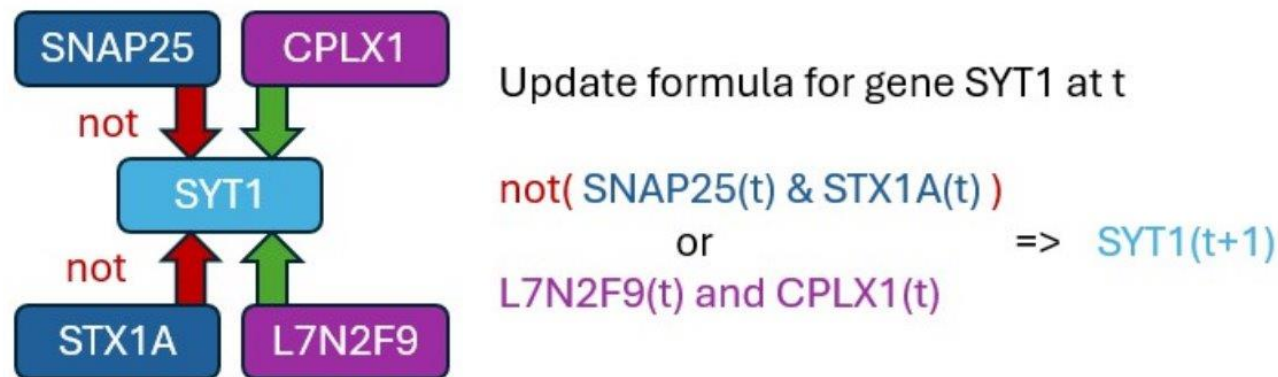
Background

1. Link with biology: dynamics in Boolean networks, canalisation-based representations
2. Symbolic music modeling
3. The SimCLR framework

Link with biology crafting general-purpose embeddings for cellular behavior predicted by Boolean networks

Remember **Boolean networks**?

Perturbation: **overexpression** (forced to True) or **knockout** (False)

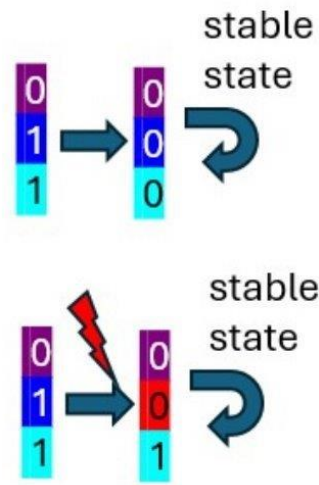
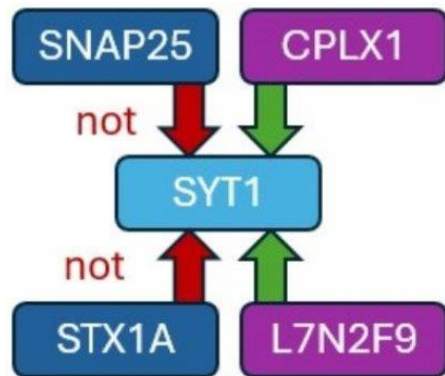


Get analytically / computationally the **stable states**

²¹ *Kauffman (1969). Journal of Theoretical Biology 22.3, pp. 437-467*
Thomas (1973). Journal of Theoretical Biology 42.3, pp. 563-585

Link with biology crafting general-purpose embeddings for cellular behavior predicted by Boolean networks

Looking at static networks VS cellular behavior

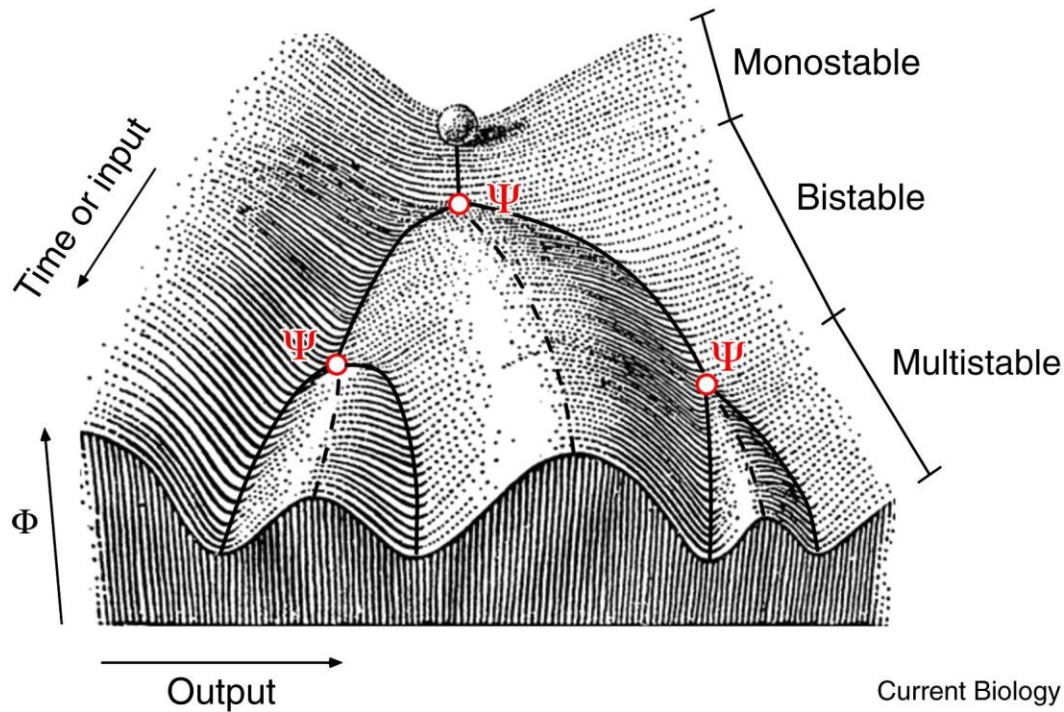


Evaluate the formula from state **0 1 1**

1. Without perturbation
 $\text{not}(\text{SNAP25}(t) \ \& \ \text{STX1A}(t))$
or
 $\text{L7N2F9}(t) \ \text{and} \ \text{CPLX1}(t)$
 $\Rightarrow \text{SYT1}(t+1)$
2. Knocking out $\text{SNAP25}(t) \ \& \ \text{STX1A}(t)$
 $\text{not}(\perp) = \text{not}(\text{False}) = \text{True}$
or
 $\text{L7N2F9}(t) \ \text{and} \ \text{CPLX1}(t)$
 $\Rightarrow \text{SYT1}(t+1)$

Canalisation crafting general-purpose embeddings for cellular behavior predicted by Boolean networks

Conrad Waddington's (epigenetic) landscape in 1942

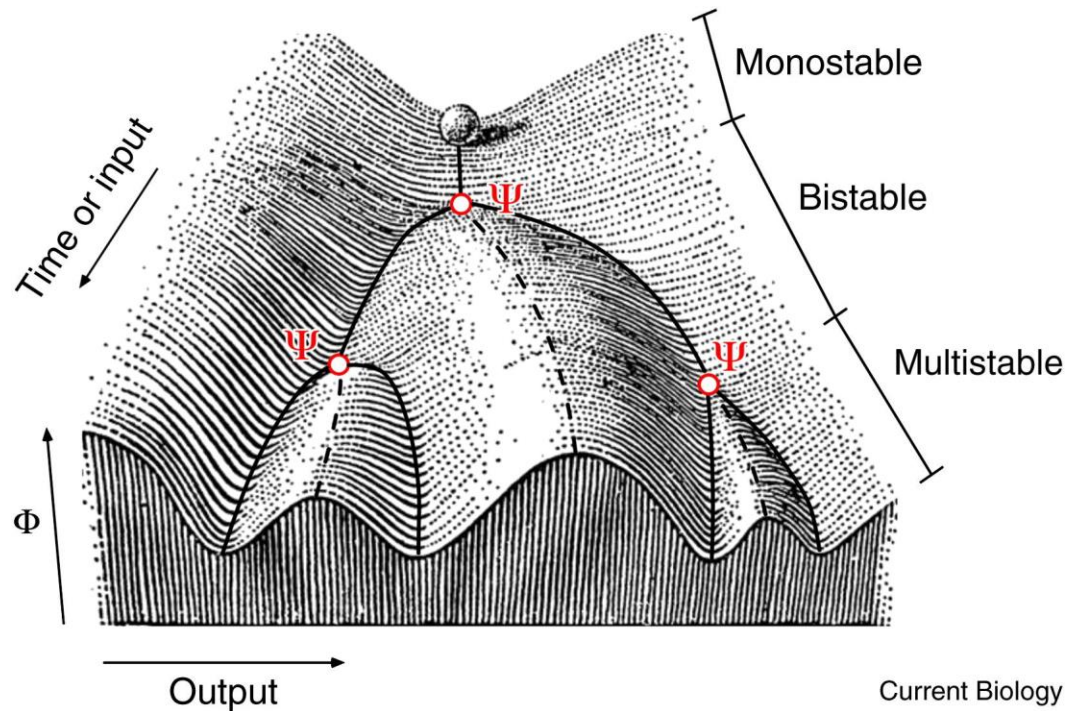


Source from [1]

[1] Ferrell. (2012). *Current biology*, 22(11), R458-R466.

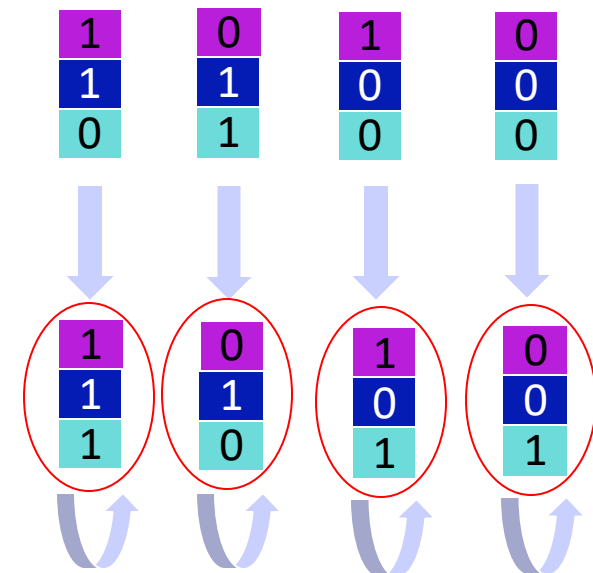
Canalisation crafting general-purpose embeddings for cellular behavior predicted by Boolean networks

Conrad Waddington's (epigenetic) landscape in 1942



Source from [1]

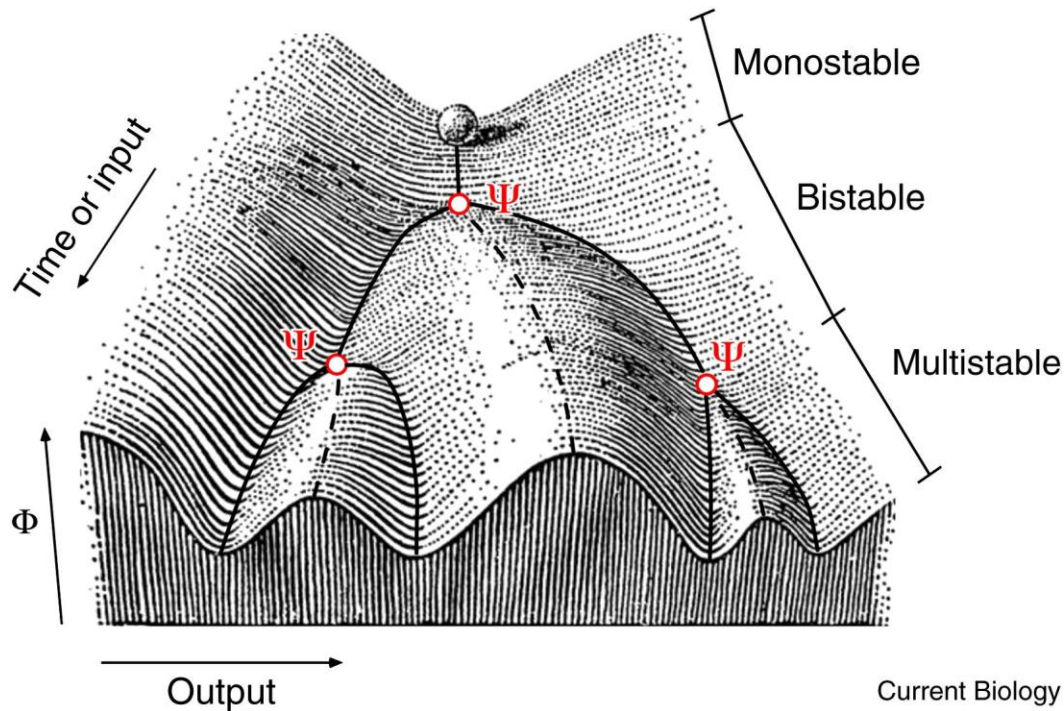
State transition diagram induced by Boolean functions



[1] Ferrell. (2012). *Current biology*, 22(11), R458-R466.

Canalisation crafting general-purpose embeddings for cellular behavior predicted by Boolean networks

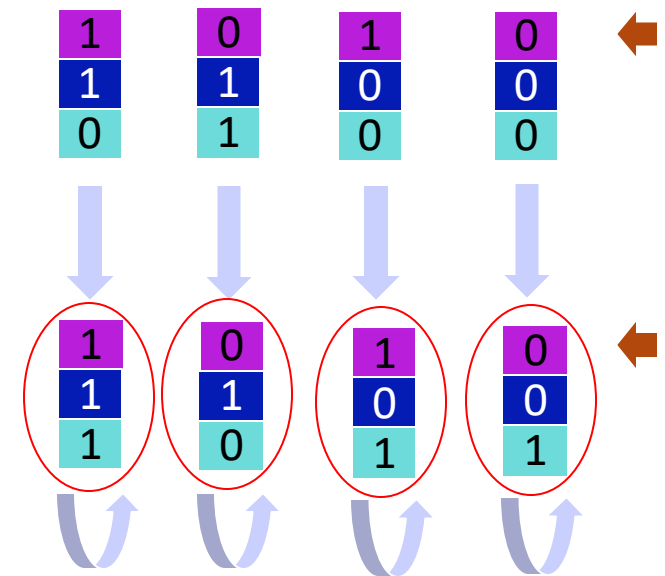
Conrad Waddington's (epigenetic) landscape in 1942



Source from [1]

Current Biology

State transition diagram induced by Boolean functions



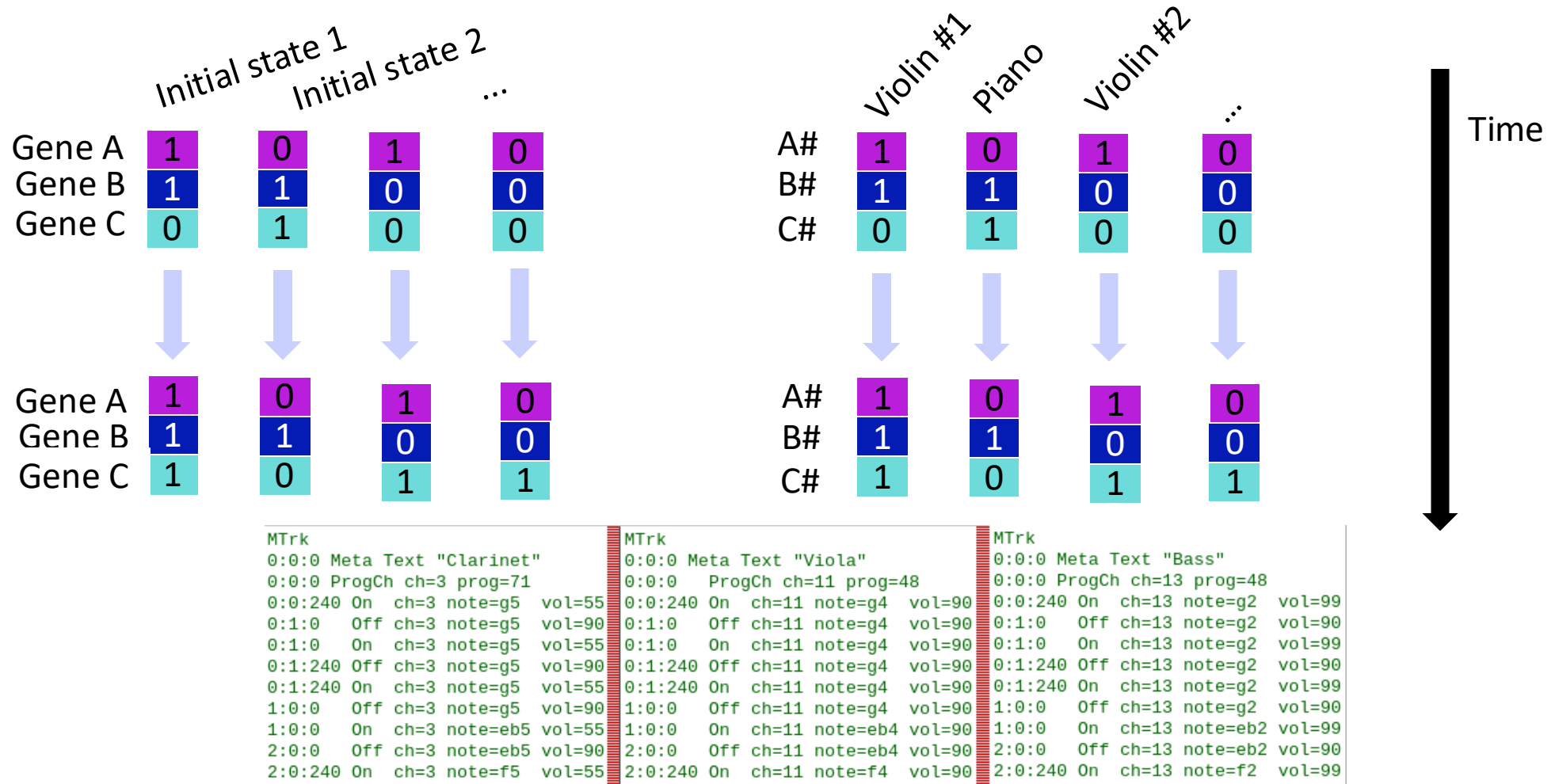
$\text{not}(\text{SNAP25}(t) \& \text{STX1A}(t))$
 or
 $\text{L7N2F9}(t) \text{ and } \text{CPLX1}(t)$

$\Rightarrow \text{SYT1}(t+1)$

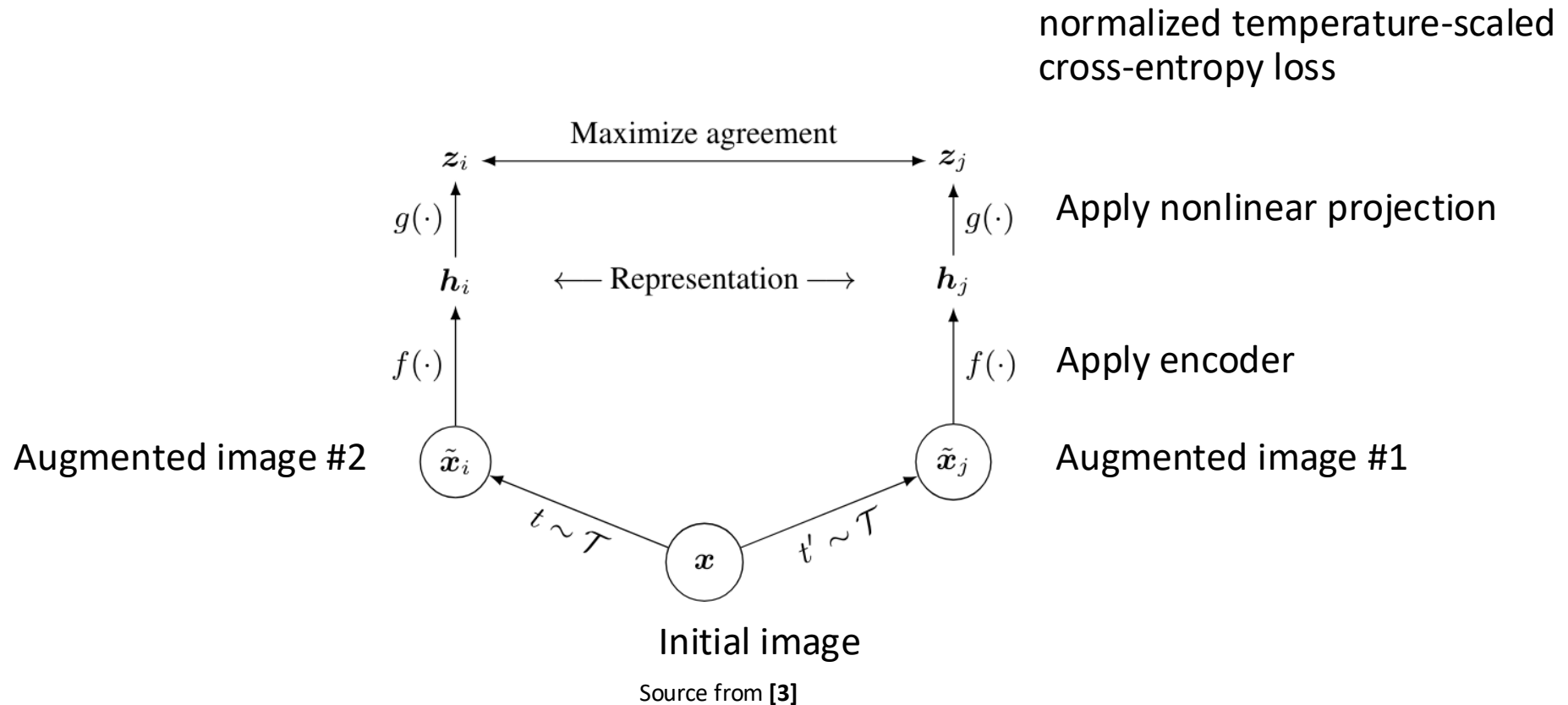
[1] Ferrell. (2012). *Current biology*, 22(11), R458-R466.

[2] He & Macauley. (2016). *Physica D: Nonlinear Phenomena*, 314, 1-8.

Symbolic music modeling Musical Instrument Digital Interface (MIDI) file format



The SimCLR framework [1] contrastive learning of visual representations



[3] Chen, et al. (2020, November). In *International conference on machine learning* (pp. 1597-1607). PMLR.

Content of the paper

"We introduce and open-source Aria, a pretrained autoregressive transformer model trained on transcriptions of solo piano recordings"

Specific issues:

- Refined contrastive learning approach on ~60k hours of music
- Generative capabilities & representation learning with few labels

Prior approaches:

- Transformers for generation, Autoencoders for representations

Content of the paper

"We introduce and open-source Aria, a pretrained autoregressive transformer model trained on transcriptions of solo piano recordings"

1. Tokenization scheme
2. Transformer architecture and preprocessing
3. Contrastive representation learning
4. Experiments: Generative capabilities
5. Experiments: Representation-based tasks

Tokenization scheme more adapted to transformers

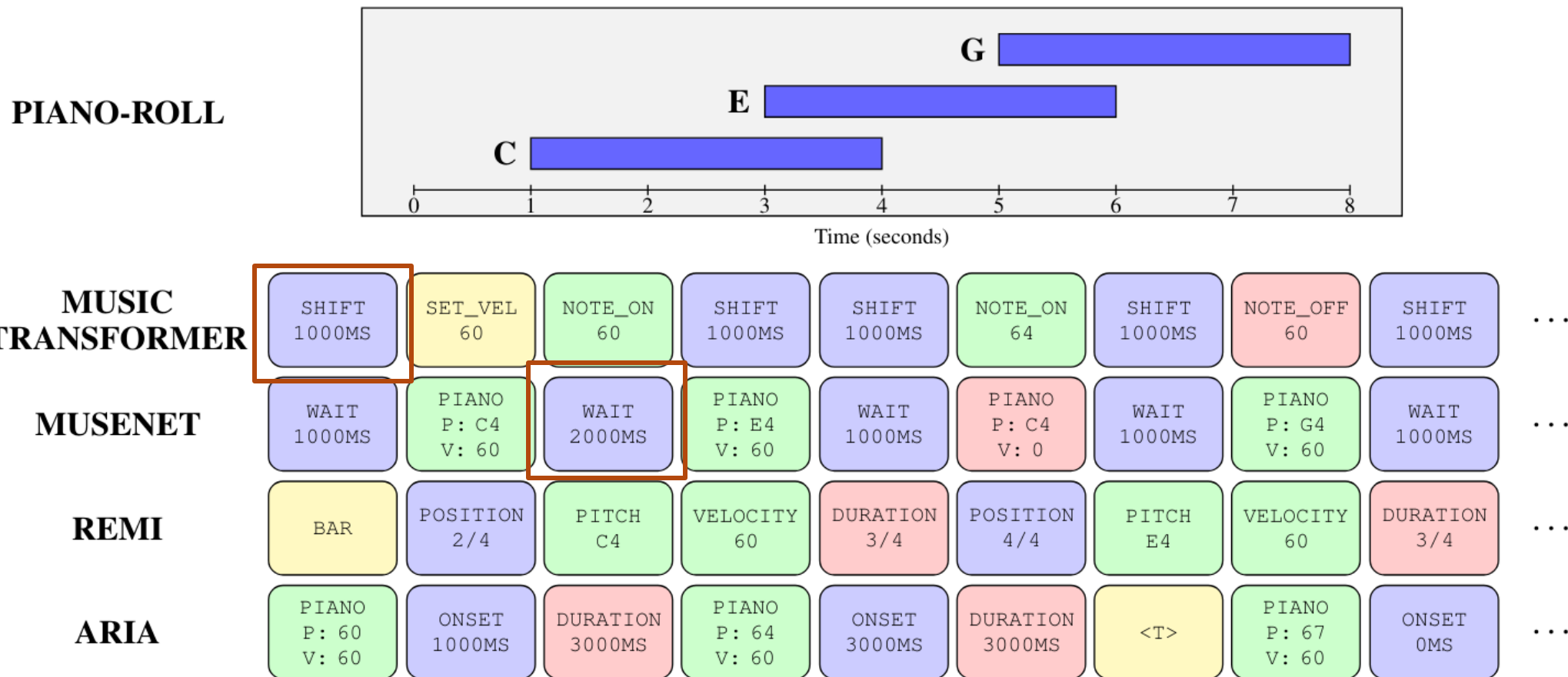


Figure 2 from paper

Tokenization scheme more adapted to transformers

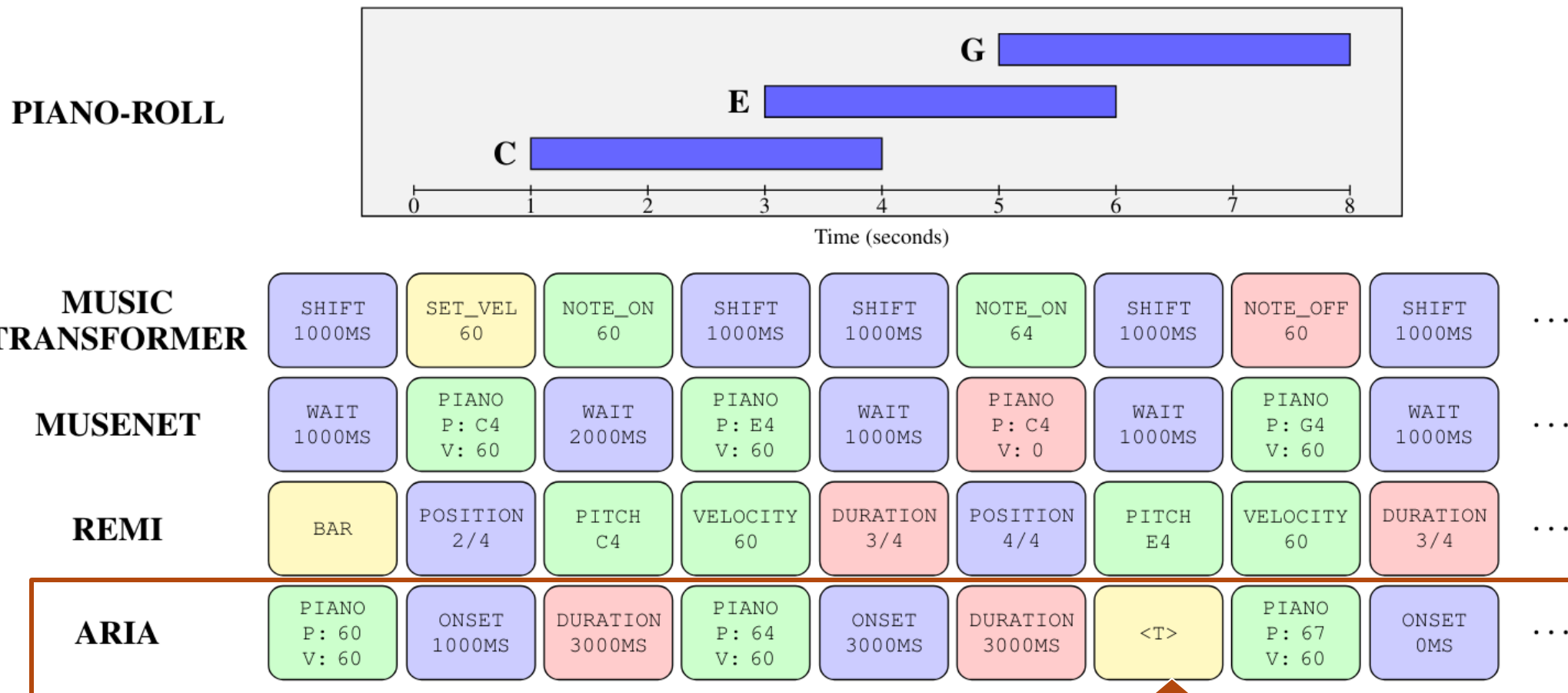


Figure 2 from paper



Transformer architecture and preprocessing

Modified LLaMa 3.2 model ~0.7B parameters, simple multi-head attention and layer normalization

Transformer architecture and preprocessing

Modified LLaMa 3.2 model ~0.7B parameters, simple multi-head attention and layer normalization

Preprocessing ARIA-MIDI data set remove duplicates, "messy" pieces

Context length: ~8k tokens

Data augmentation: random transposition, varying tempo and velocity

Transformer architecture and preprocessing

Modified LLaMa 3.2 model ~0.7B parameters, simple multi-head attention and layer normalization

Preprocessing ARIA-MIDI data set remove duplicates, "messy" pieces

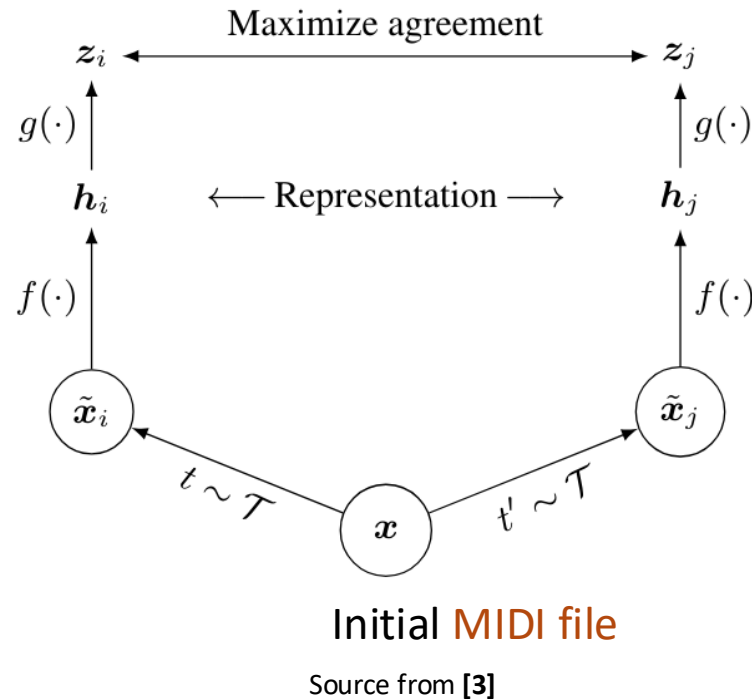
Context length: ~8k tokens

Data augmentation: *random transposition, varying tempo and velocity*

Finetuning for generation single epoch on short solo piano prompts (LoRA?)

Add a new token <D> 100 tokens before the end of the prompt

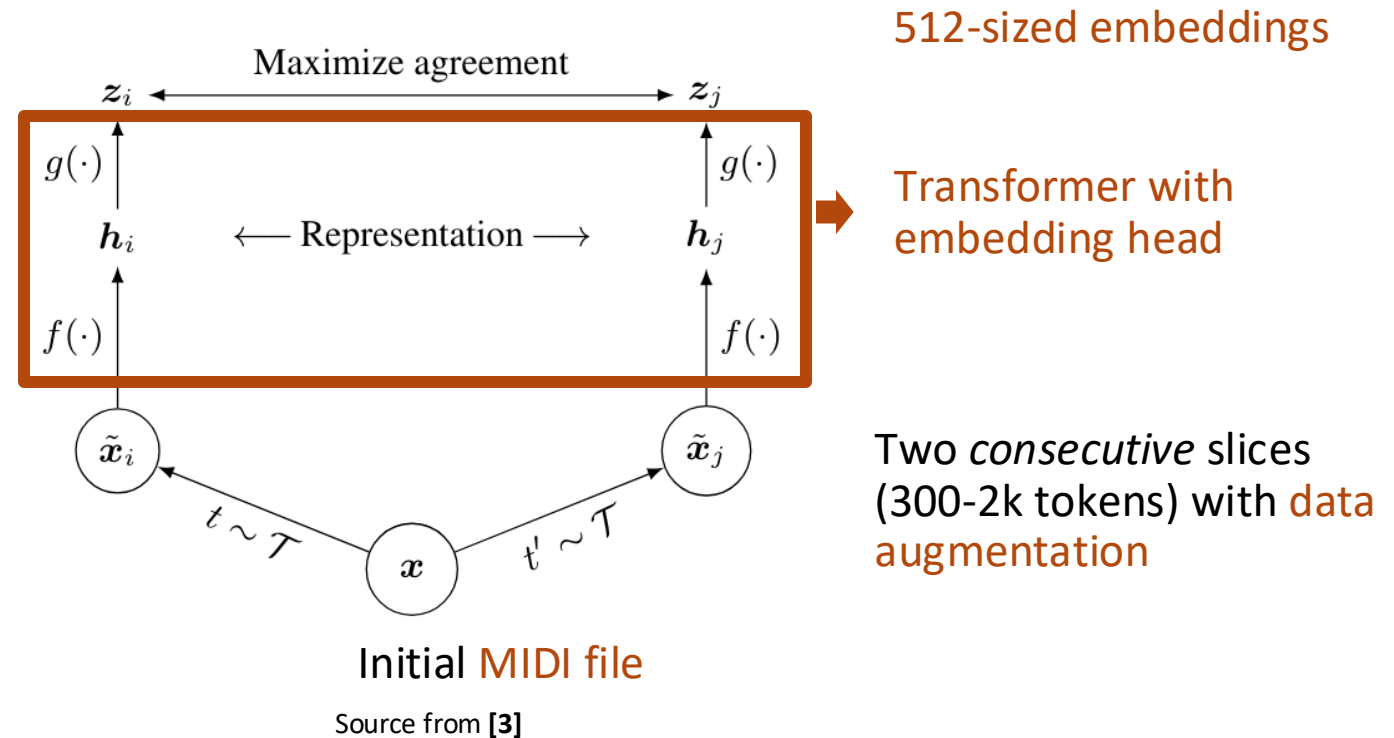
Contrastive representation learning adapted from the SimCLR framework



Two *consecutive* slices
(300-2k tokens) with **data
augmentation**

[3] Chen, et al. (2020, November). In *International conference on machine learning* (pp. 1597-1607). PMLR.

Contrastive representation learning adapted from the SimCLR framework



[3] Chen, et al. (2020, November). In *International conference on machine learning* (pp. 1597-1607). PMLR.

Experiments: Generative capabilities

45" continuation of 15" prompts random pairwise A/B human evaluation

Compared Model	Wins	Ties	Losses	p-value
AM Transformer	38	0	6	$9.43e-7$
Suno 3.5	18	9	21	$7.49e-1$
MusicGen	49	1	0	$3.55e-15$
Ground Truth	15	9	17	$8.60e-1$

They say there are 46 participants, yet most lines do not sum to 46...

40 continuations /model (5 subgenres x 8 prompts)

Experiments: Representation-based tasks

"Composer" classification task averaging slice embeddings within a MIDI file

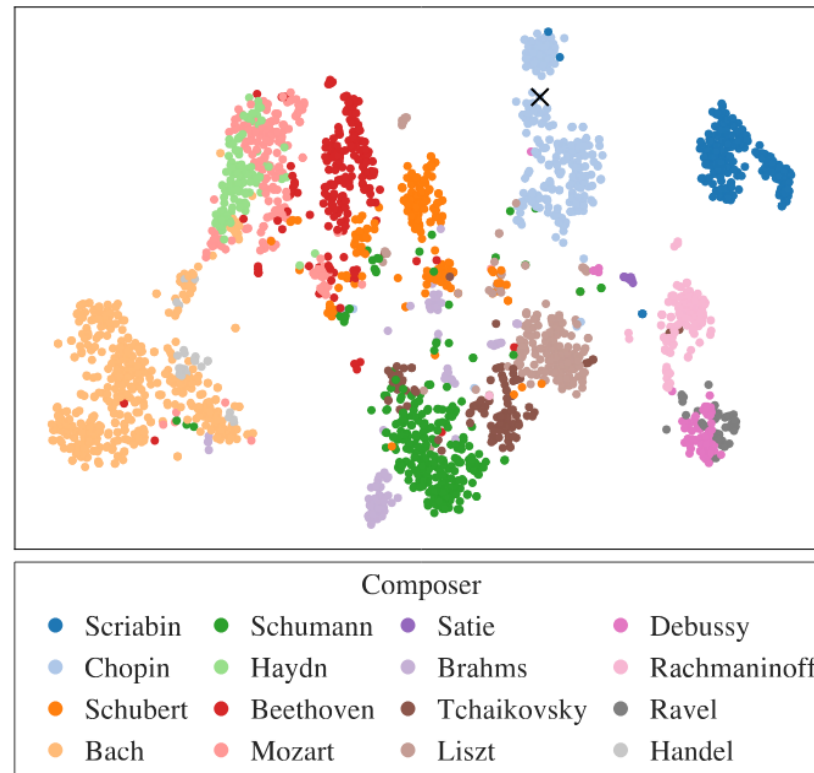


Figure 1 from paper

Experiments: Representation-based tasks

Model	custom classification tasks								benchmarks			
	Genre		Form		Musical Period		Composer		Pianist8		VG-MIDI	
	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1	Acc	F1
<i>Main Results</i>												
MERT	83.00	83.00	63.89	63.90	69.50	68.94	69.60	69.30	65.06	65.18	45.45	40.37
M3	85.10	85.10	69.88	70.12	71.20	70.81	71.90	71.72	81.93	81.48	54.55	46.13
CLaMP 3	89.10	89.10	77.79	77.97	80.60	80.20	84.50	84.46	80.72	79.76	45.45	36.53
Aria _{Emb}	<u>92.40</u>	<u>92.40</u>	<u>82.45</u>	<u>82.57</u>	<u>84.70</u>	<u>84.69</u>	<u>90.50</u>	<u>90.49</u>	91.57	92.38	<u>63.64</u>	<u>63.96</u>
Aria _{Ft}	93.20	93.20	87.53	87.59	86.50	86.53	96.30	96.32	<u>91.56</u>	<u>92.03</u>	68.18	69.55



supervised finetuned model with the replacement by a classifier head

Perspectives

1. Comments on the paper
2. Why is it interesting for BioComp?

My comments on the paper

Strengths:

- The method is elegant and adapted to multi-tasking
- Results on embedding learning are interesting
- Open-source: <https://github.com/EleutherAI/aria>

Weaknesses:

- Results on generative capabilities seem strange to me
- **Only applied to solo piano transcriptions (no multi-track)**
- Could other data augmentations be applied? (e.g., changes to inter-onset intervals)

Your comments?

Why is it interesting for BioComp?

Adaptation of methods applied to imaging to other time-series data

- Subtle dynamical changes
- Required consistency across time (harmonics...)
- Time integration into transformer models
- Restricted number of samples

Under review as a conference paper at ICLR 2026

PIANIST TRANSFORMER: TOWARDS EXPRESSIVE
PIANO PERFORMANCE RENDERING VIA SCALABLE
SELF-SUPERVISED PRE-TRAINING

Anonymous authors
Paper under double-blind review

ABSTRACT

Existing methods for expressive music performance rendering rely on supervised learning over small labeled datasets, which limits scaling of both data volume and model size, despite the availability of vast unlabeled music, as in vision

[4] <https://openreview.net/pdf?id=Kq9MMEynGW>

(perhaps a bit of a stretch for Boolean networks)