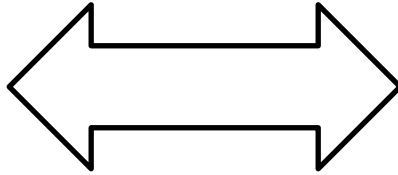


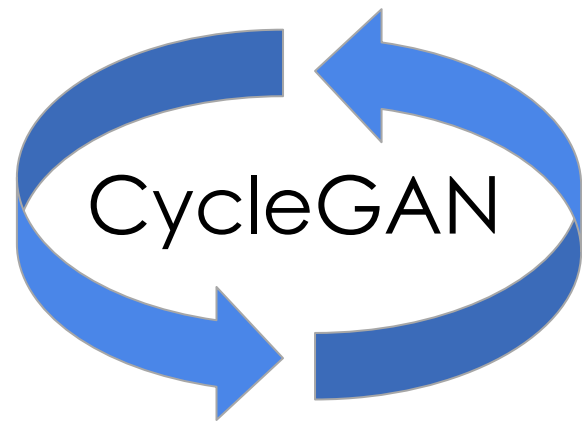
Symbolic Music Genre Transfer Insights

*Gino Brunner, Mazda Moayeri, **Oliver Richter**, Roger Wattenhofer, Chi Zhang**

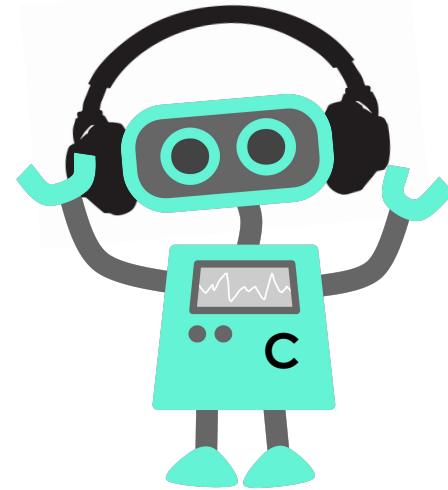
**Alphabetical order*

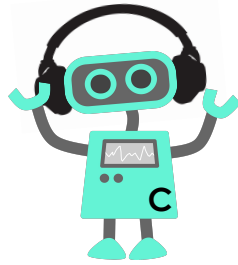
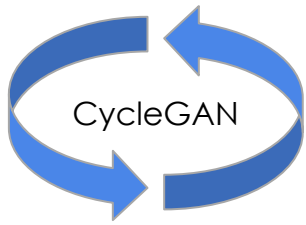


Jazz, Classic, Pop...?



Genre Classifier



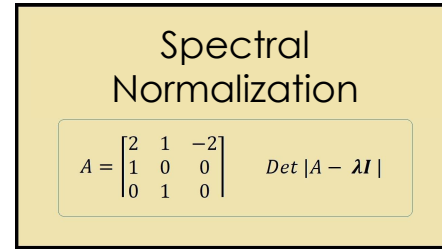
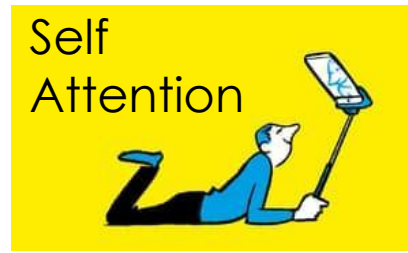
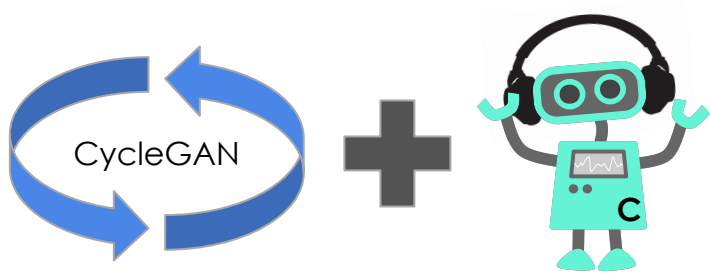


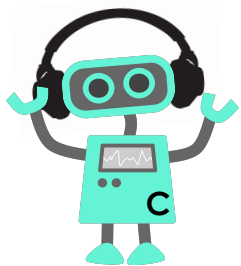
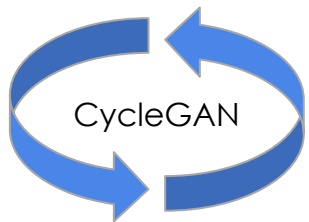
Self
Attention



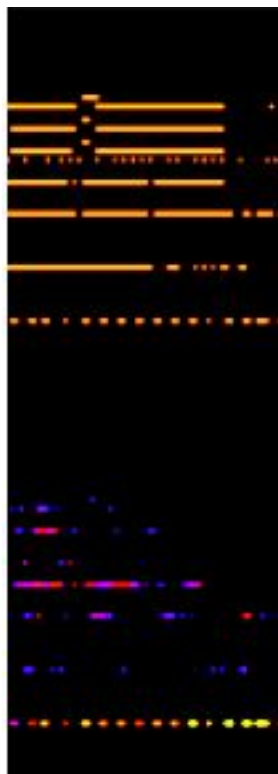
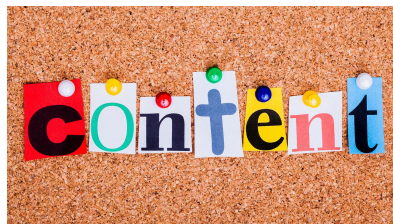
Spectral
Normalization

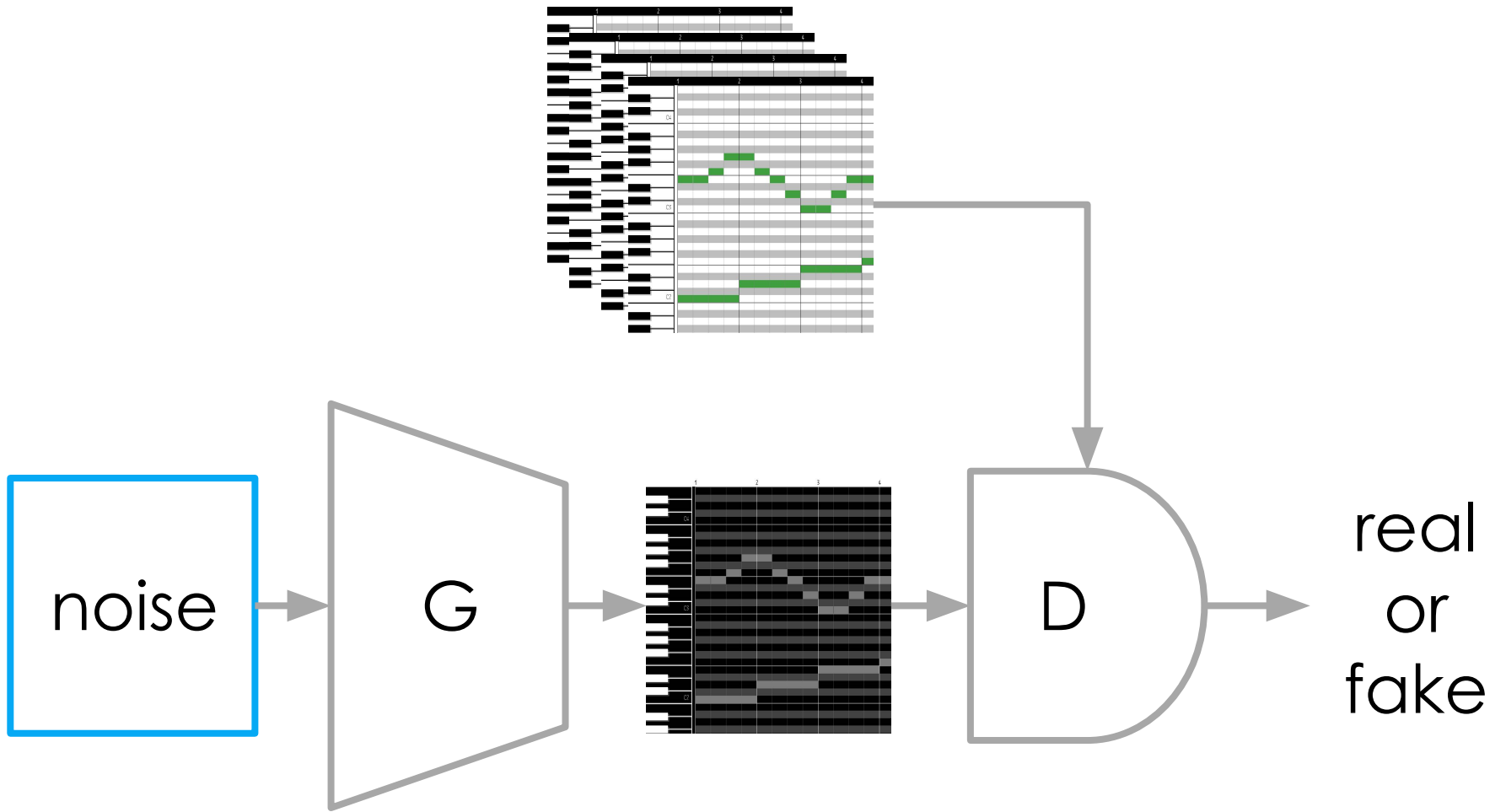
$$A = \begin{bmatrix} 2 & 1 & -2 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \quad \text{Det } |A - \lambda I|$$

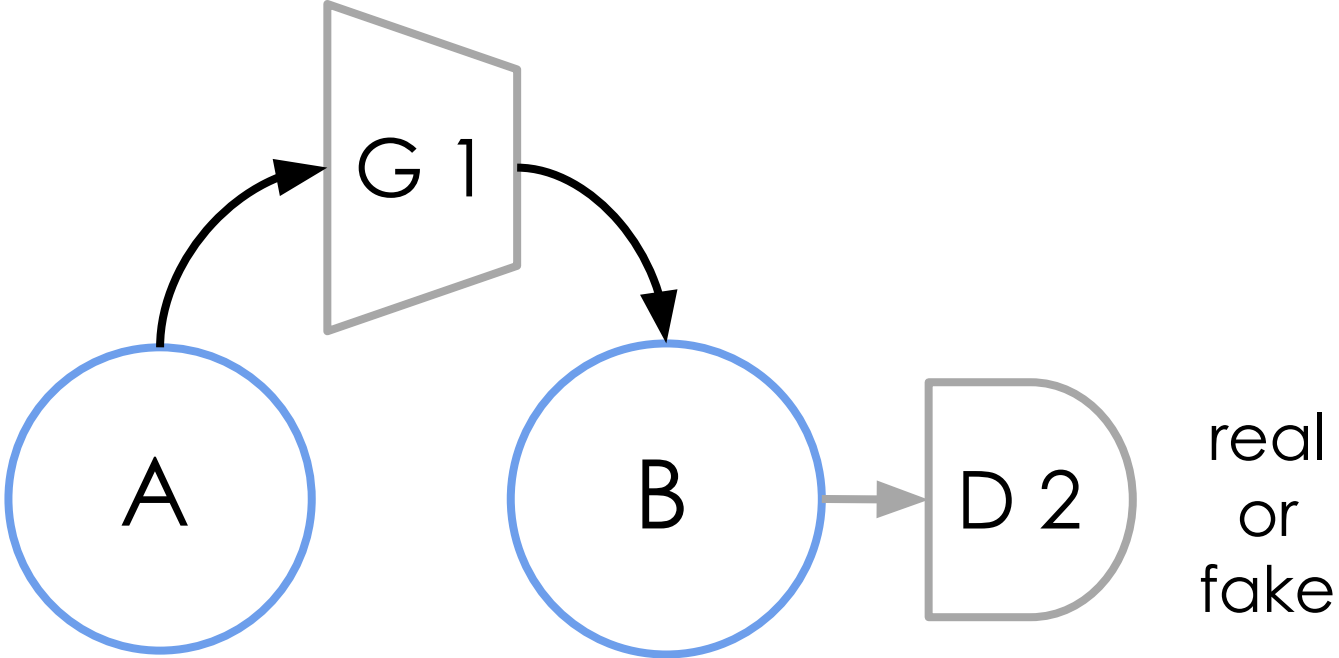


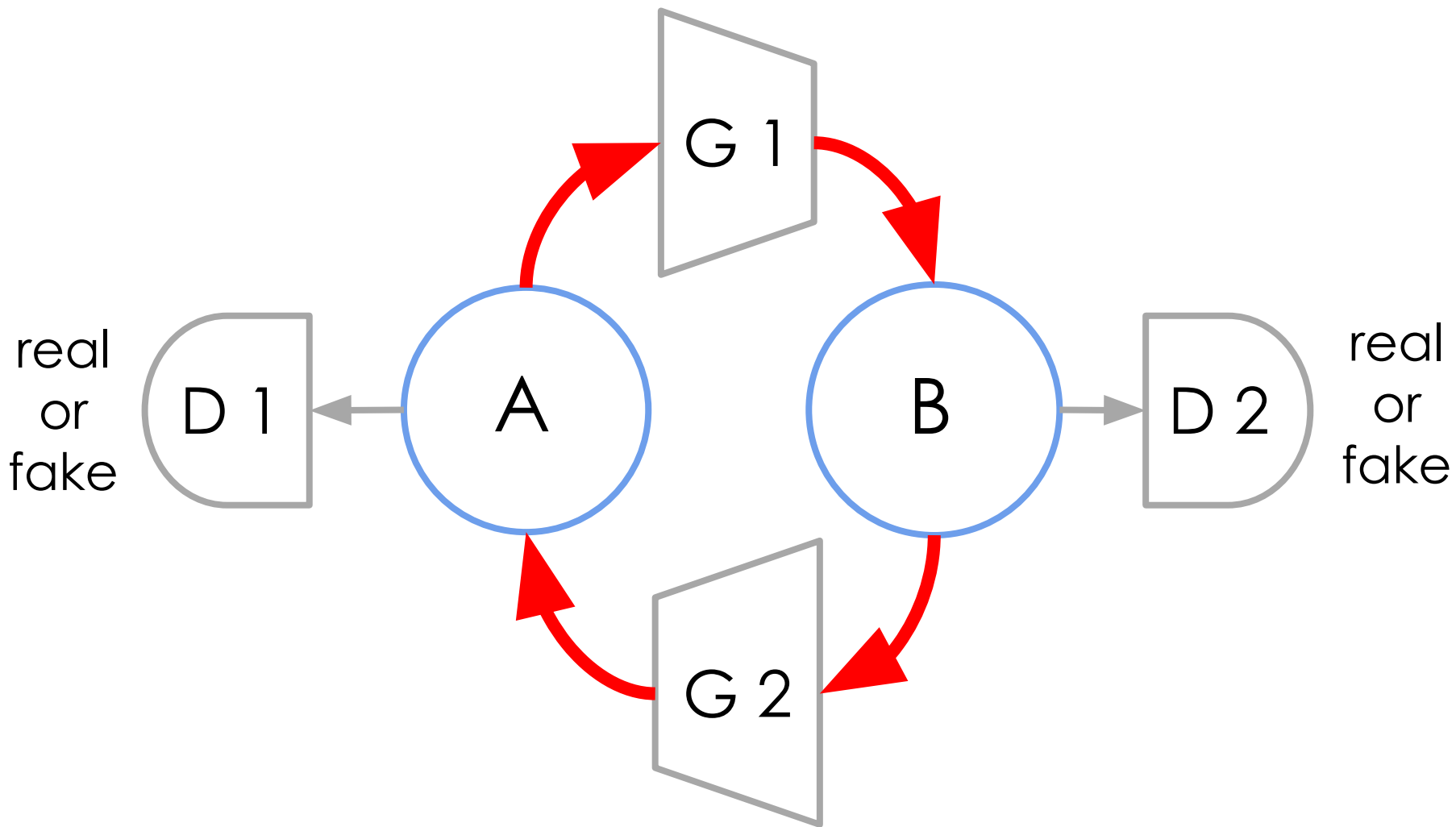


Spectral Normalization

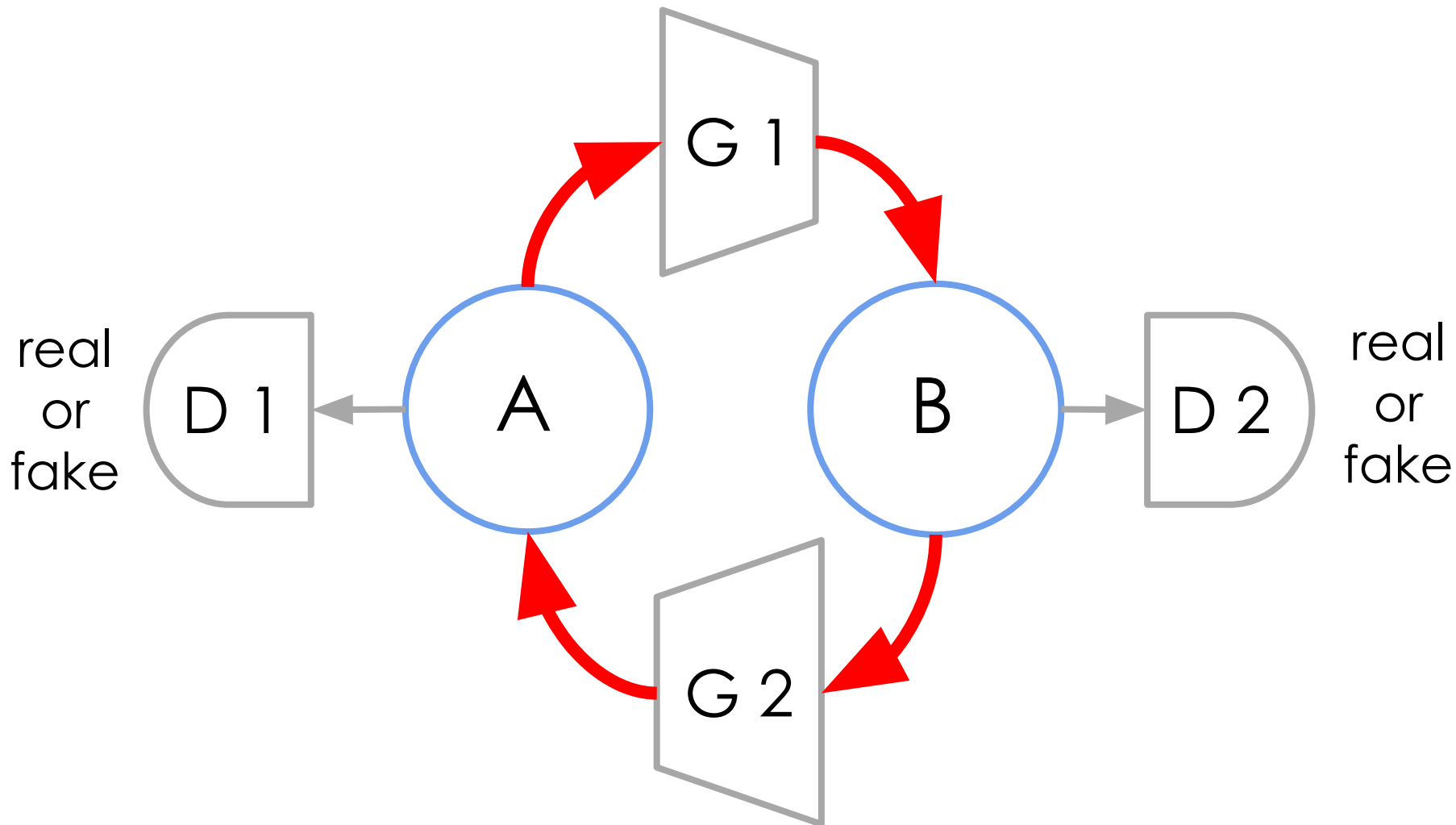
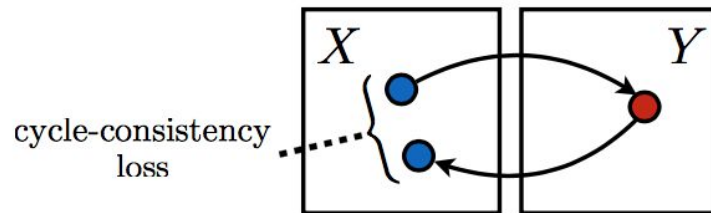
$$A = \begin{bmatrix} 2 & 1 & -2 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \quad \text{Det } |A - \lambda I|$$


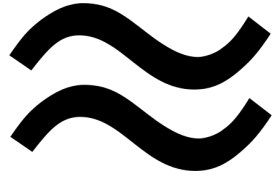






Cycle Consistency!

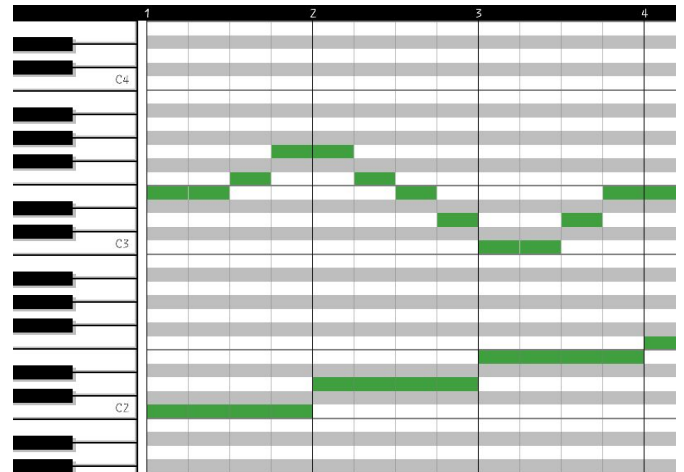




MIKI



Piano Roll



Classical Acoustic Guitar | Ch1

Roots Rock | Ch1

Roots Rock | Ch1

Classical Acoustic Guitar | Ch1

Upright Studio Bass | Ch1

Steinway Grand Piano | Ch1

Nylon Gtr

Jazz Gtr

Jazz Gtr

Nylon Gtr

Acoustic Bass (+12 semitones)

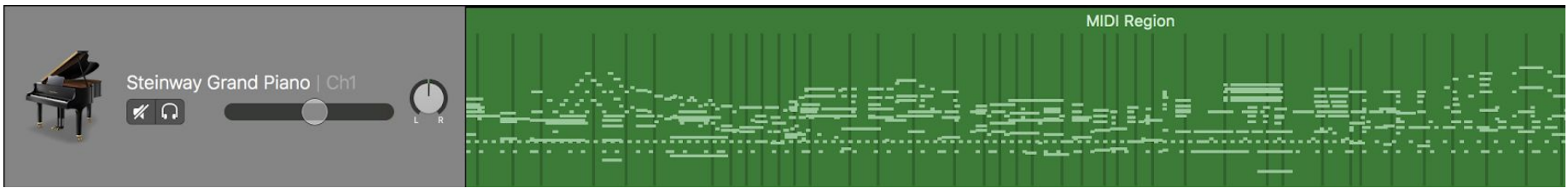
MIDI Region



Discard Drums
and Merge

Steinway Grand Piano | Ch1

MIDI Region



Discard velocity



Piano Roll

84 pitches

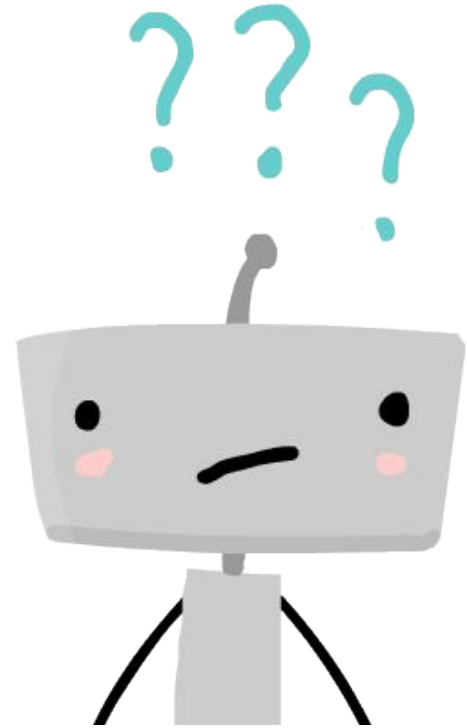


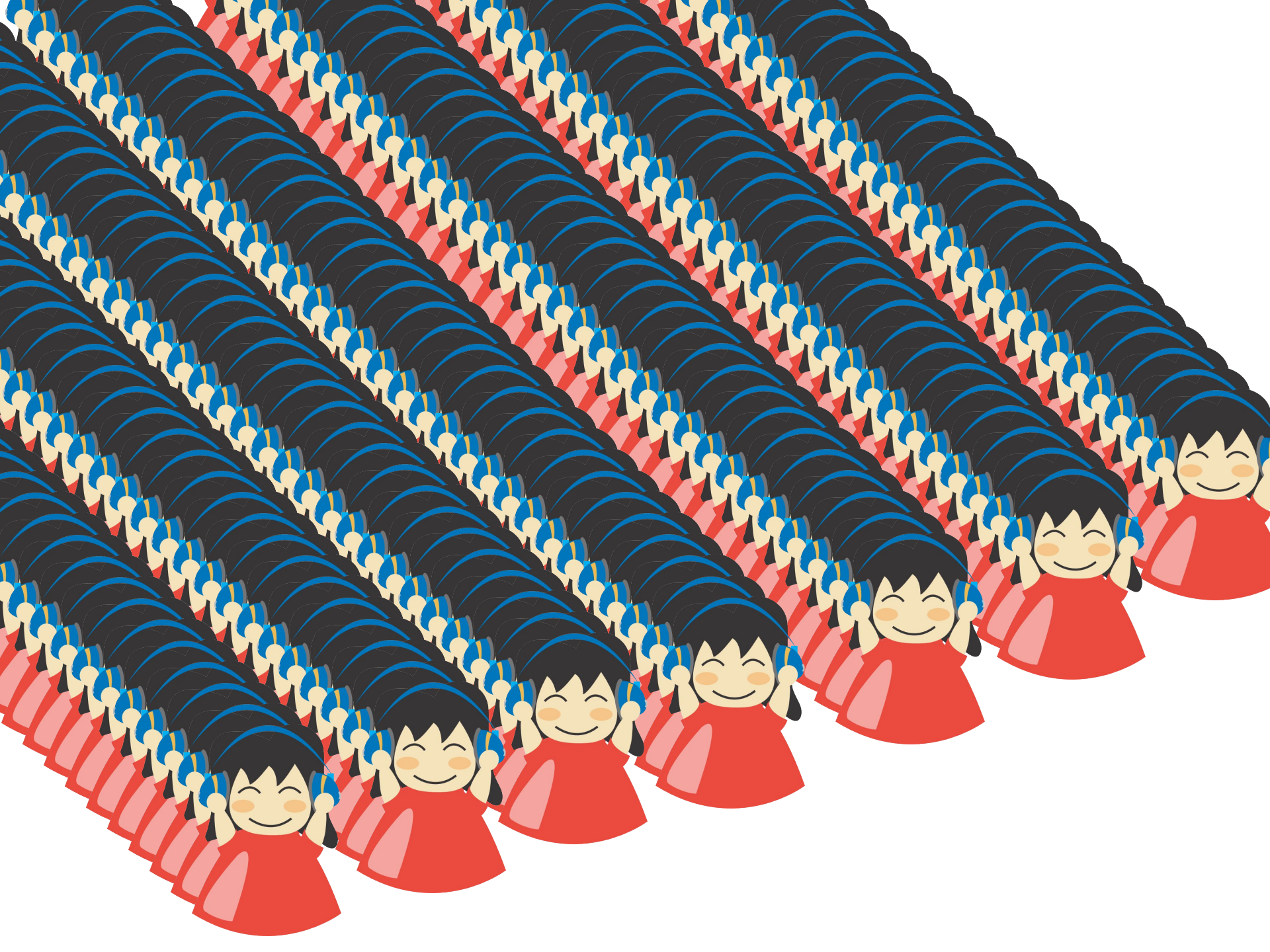
64 beats = 4 bars = 1 phrase



	Jazz	Classic	Pop
Songs	559	2714	1069
Phrases	12341	16545	20780

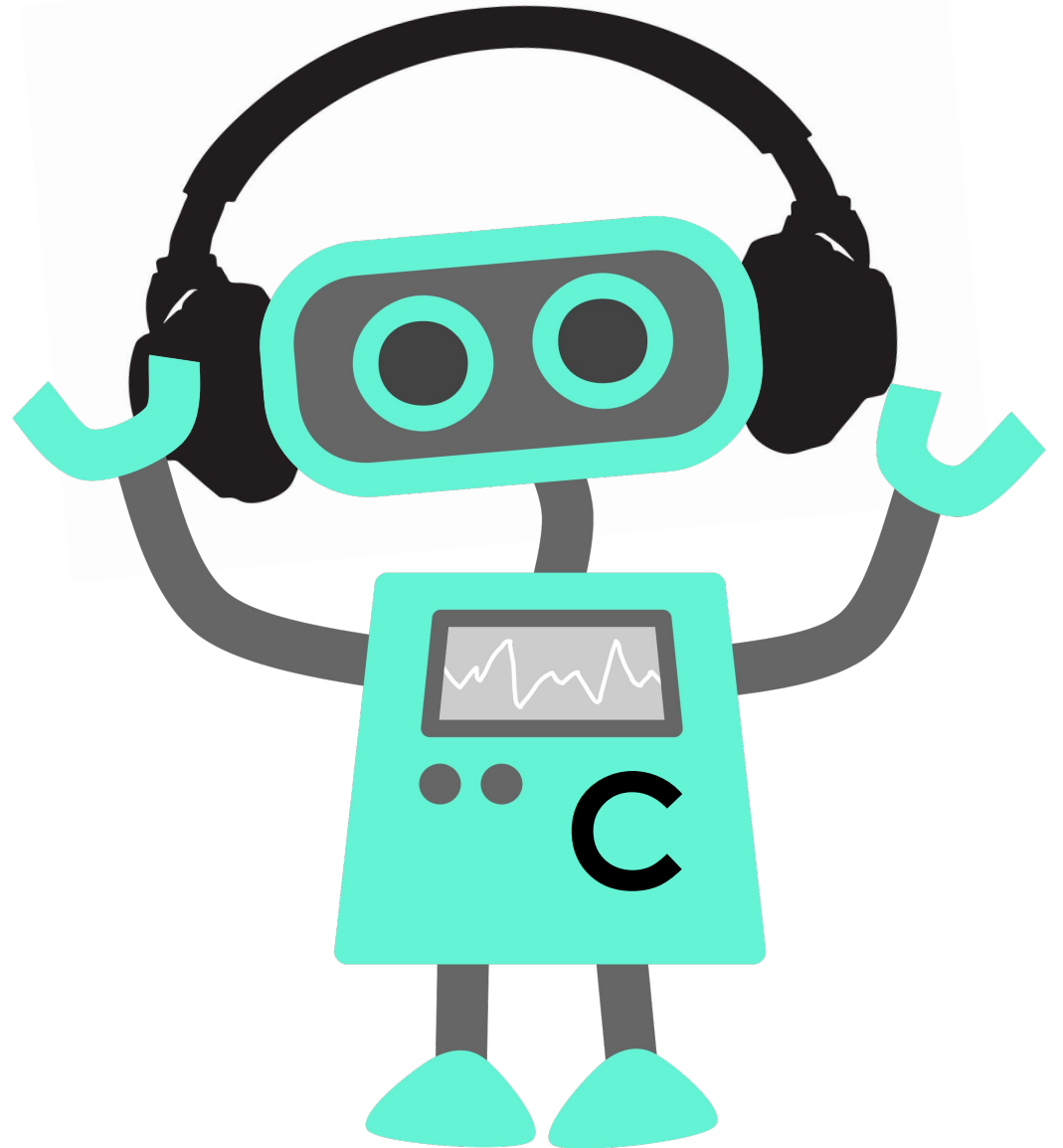


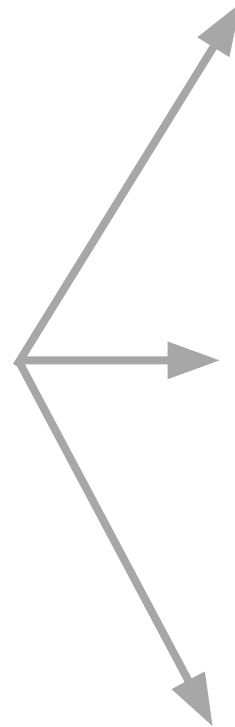
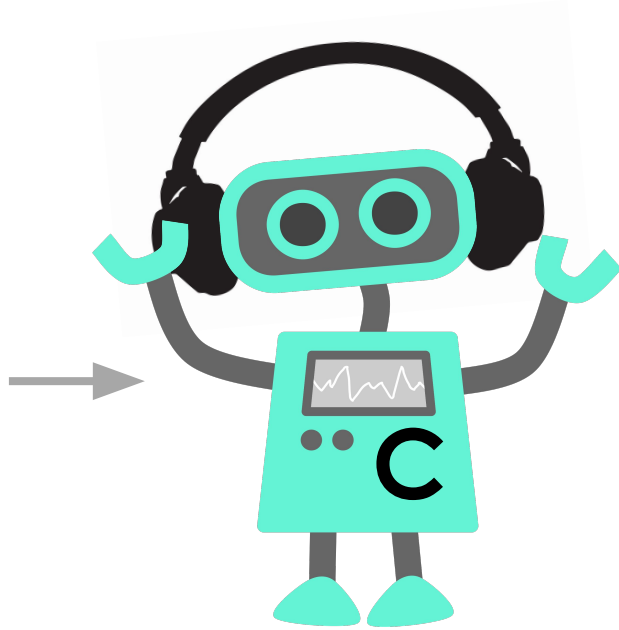
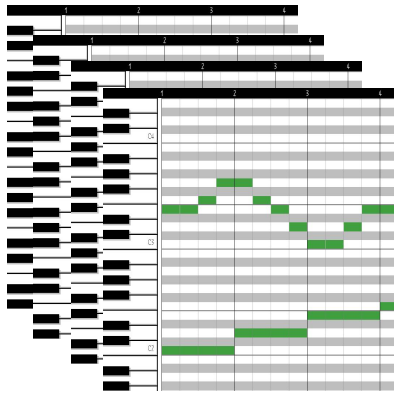




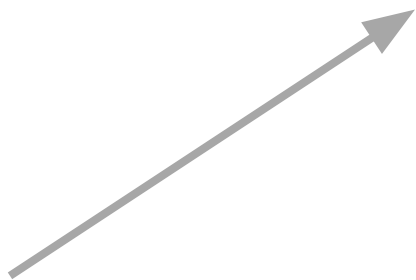
Genre Classifier

Let me do it!

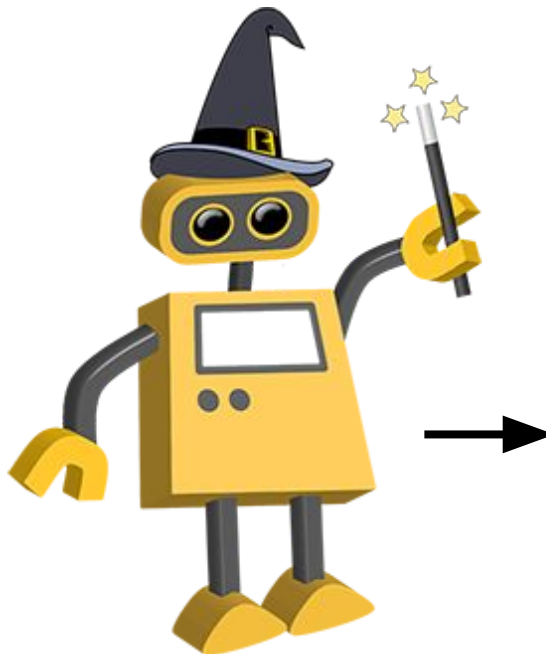




90% Jazz!



Before

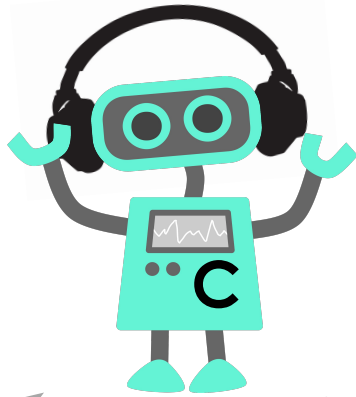


After



90% Jazz!

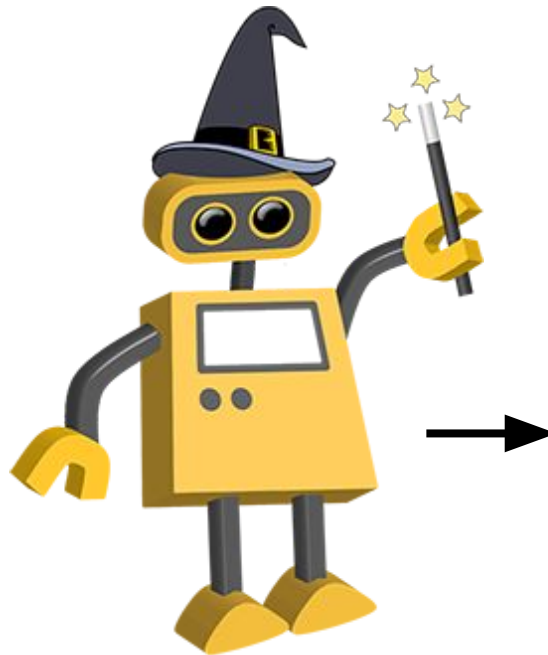
80% Classic!



Before



After



1. classic: 86%



2. classic -> pop: 16%



3. classic -> pop -> classic: 80%

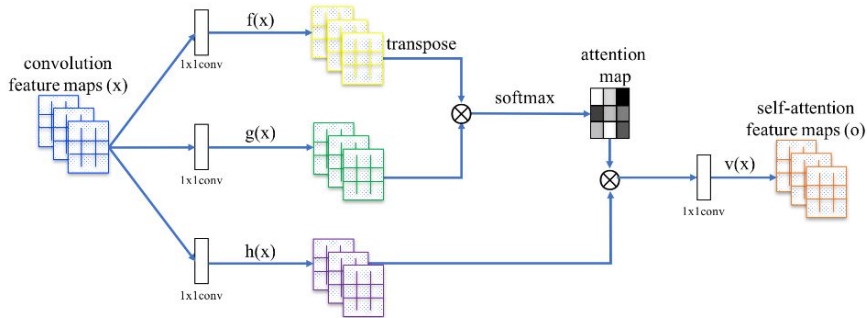


$$s_1 = \frac{1}{2} ((86\% - 16\%) + (80\% - 16\%)) = 67\%$$

$$s_2 = 60\%$$

$$s_{tot}^D = \frac{1}{2} (s_1 + s_2) = 63.5\%$$

Self-Attention



Spectral Normalization

Algorithm 1 SGD with spectral normalization

- Initialize $\tilde{\mathbf{u}}_l \in \mathcal{R}^{d_l}$ for $l = 1, \dots, L$ with a random vector (sampled from isotropic distribution).
- For each update and each layer l :

1. Apply power iteration method to a unnormalized weight W^l :

$$\tilde{\mathbf{v}}_l \leftarrow (W^l)^T \tilde{\mathbf{u}}_l / \|(W^l)^T \tilde{\mathbf{u}}_l\|_2 \quad (20)$$

$$\tilde{\mathbf{u}}_l \leftarrow W^l \tilde{\mathbf{v}}_l / \|W^l \tilde{\mathbf{v}}_l\|_2 \quad (21)$$

2. Calculate \bar{W}_{SN}^l with the spectral norm:

$$\bar{W}_{\text{SN}}^l(W^l) = W^l / \sigma(W^l), \text{ where } \sigma(W^l) = \tilde{\mathbf{u}}_l^T W^l \tilde{\mathbf{v}}_l \quad (22)$$

3. Update W^l with SGD on mini-batch dataset \mathcal{D}_M with a learning rate α :

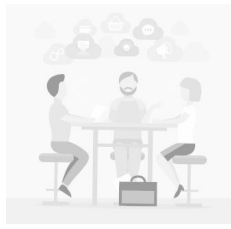
$$W^l \leftarrow W^l - \alpha \nabla_{W^l} \ell(\bar{W}_{\text{SN}}^l(W^l), \mathcal{D}_M) \quad (23)$$

	J vs. P	C vs. P	J vs. C
Baseline	28.49%	64.62%	57.64%
SN	32.16%	61.88%	63.98%
SA	44.85%	59.35%	63.56%
SN+SA	33.23%	53.07%	66.76%

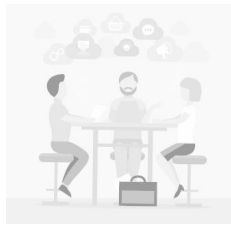
	J vs. P	C vs. P	J vs. C
Baseline	28.49%	64.62%	57.64%
SN	32.16%	61.88%	63.98%
SA	44.85%	59.35%	63.56%
SN+SA	33.23%	53.07%	66.76%

Need to train more
models!





Black	White	White	White	White	White	Black
White	Black	White	Black	Black	Black	White
White	White	Black	White	White	White	White
White	White	White	White	Black	White	White
Black	White	Black	Black	White	White	White



3

added



3

removed



6

total



1



2



3

Baseline

SN

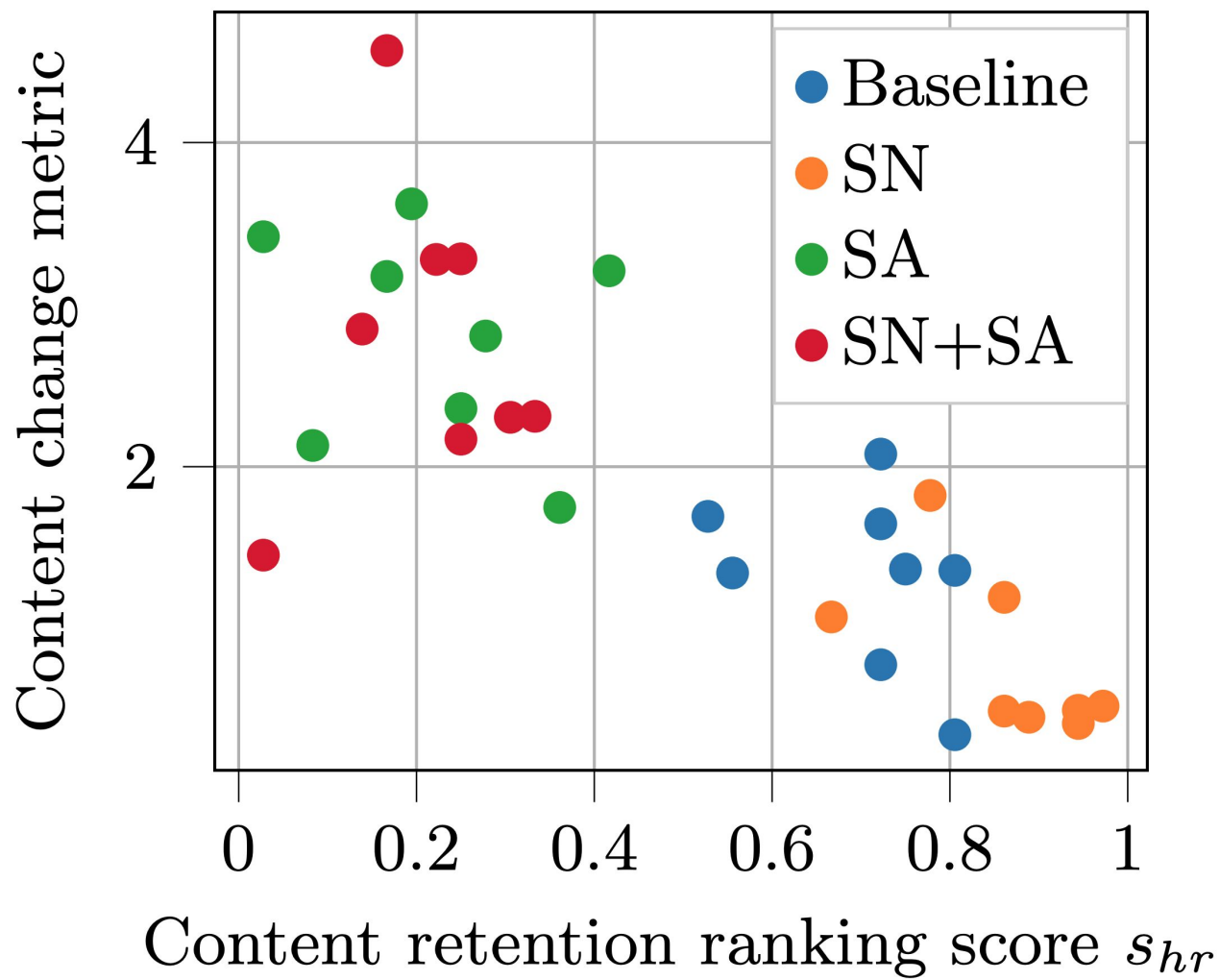
SA

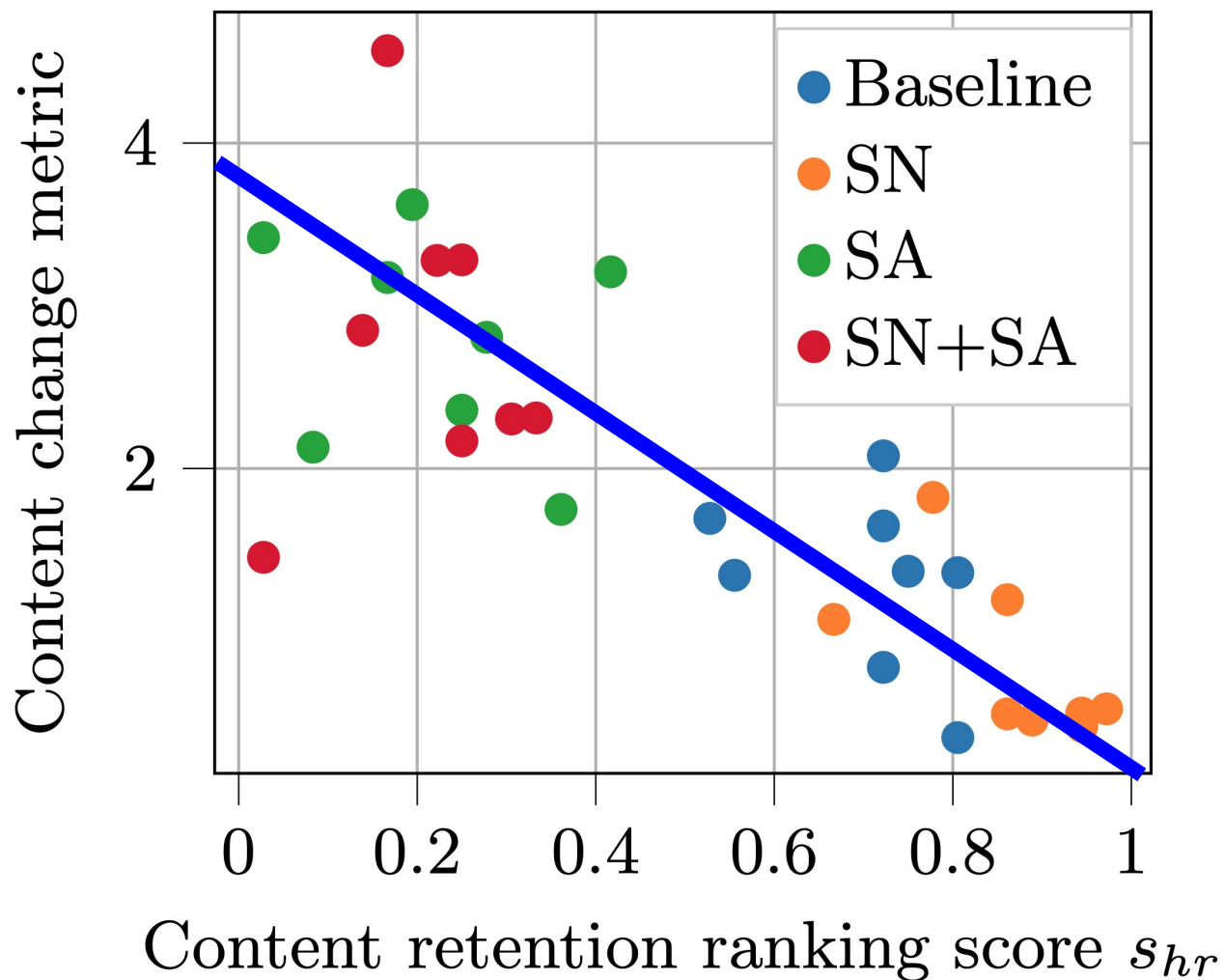
SN+SA



Recognizable?

Fidelity?

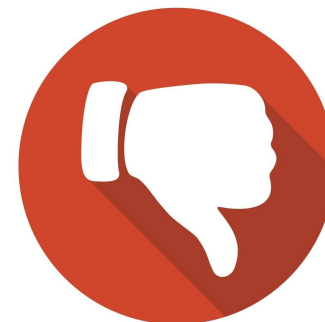
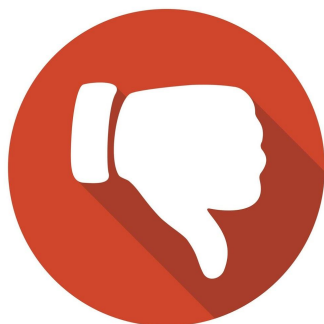




Pearson correlation: **-0.805**



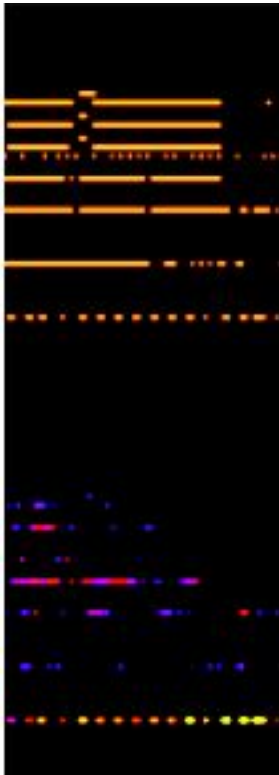
	Baseline		SN	
	$C \rightarrow P$	$P \rightarrow C$	$C \rightarrow P$	$P \rightarrow C$
added	0.82 ± 0.45	0.28 ± 0.17	0.57 ± 0.38	0.08 ± 0.09
removed	0.27 ± 0.17	0.46 ± 0.12	0.91 ± 0.07	0.3 ± 0.11
total	1.10 ± 0.5	0.75 ± 0.25	0.66 ± 0.38	0.38 ± 0.16
	SA		SN + SA	
added	1.47 ± 0.41	0.85 ± 0.79	1.78 ± 1.33	0.72 ± 0.66
removed	0.95 ± 0.04	0.95 ± 0.04	0.95 ± 0.05	0.93 ± 0.05
total	2.42 ± 0.42	1.79 ± 0.78	2.73 ± 1.33	1.65 ± 0.66



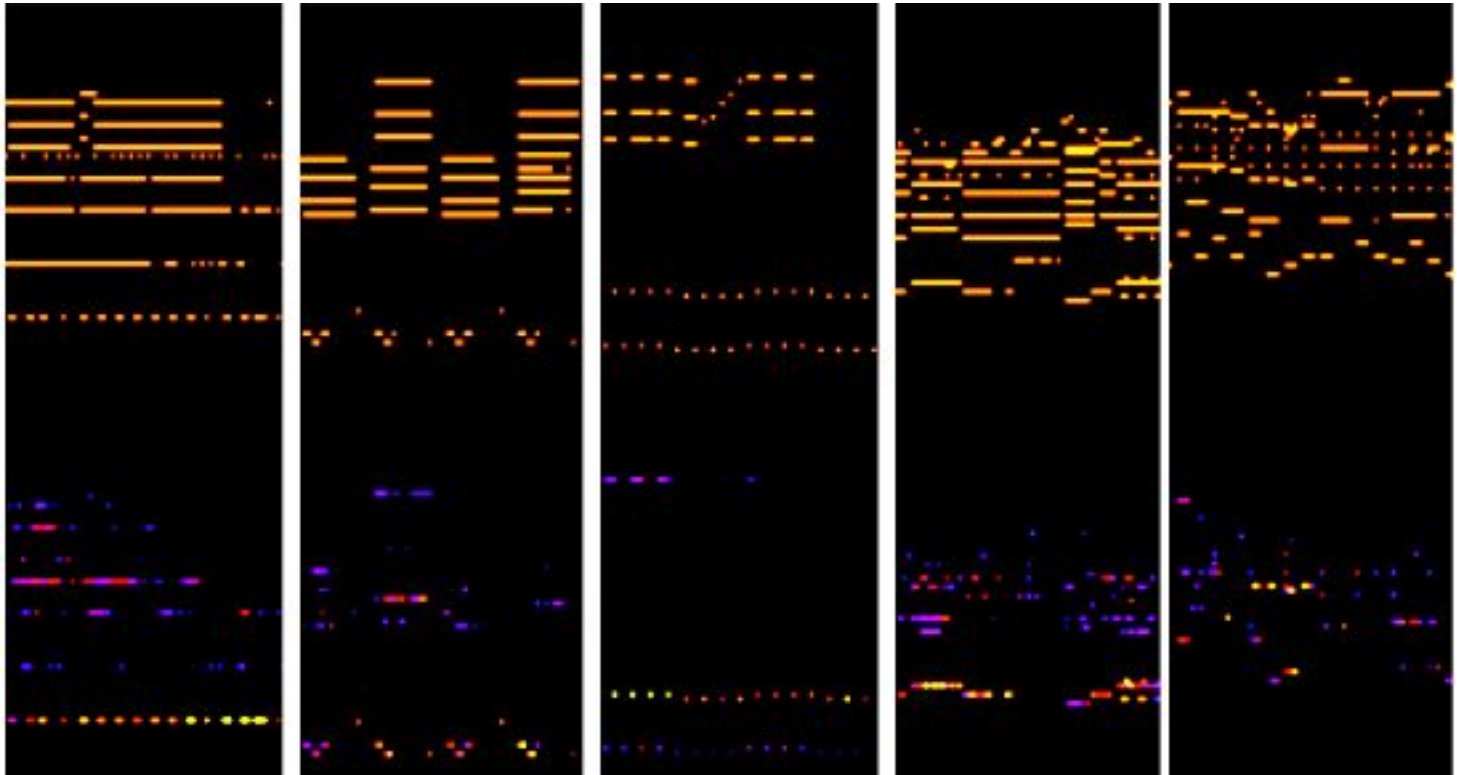


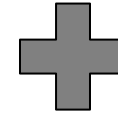
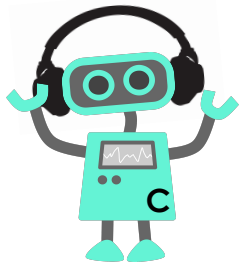
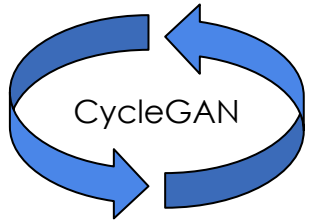
1. SN	0.6
2. SN+SA	0.51
3. Baseline	0.47
4. SA	0.42

1.0: Always ranked best
0.0: Always ranked worst

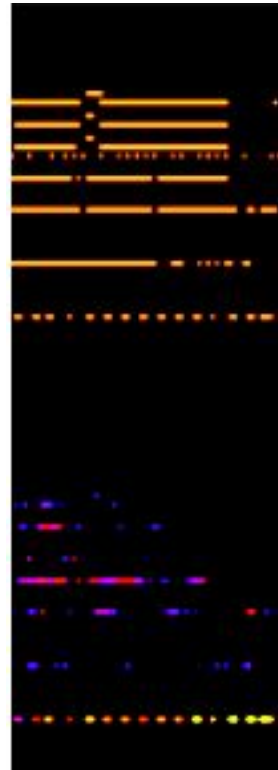
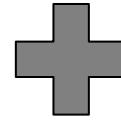


Jazz Samples

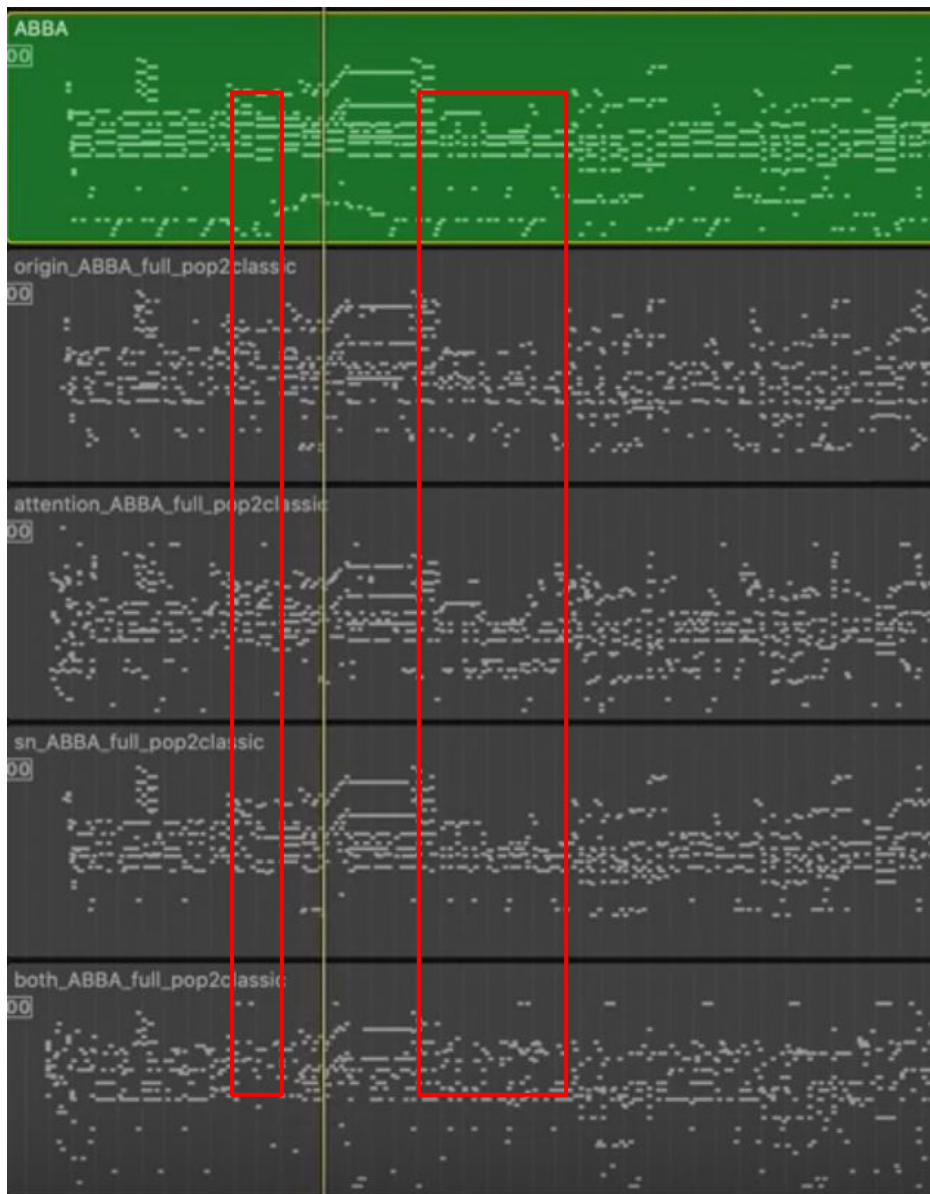




Spectral Normalization

$$A = \begin{bmatrix} 2 & 1 & -2 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \quad \text{Det } |A - \lambda I|$$


Backup Slides



original

Baseline

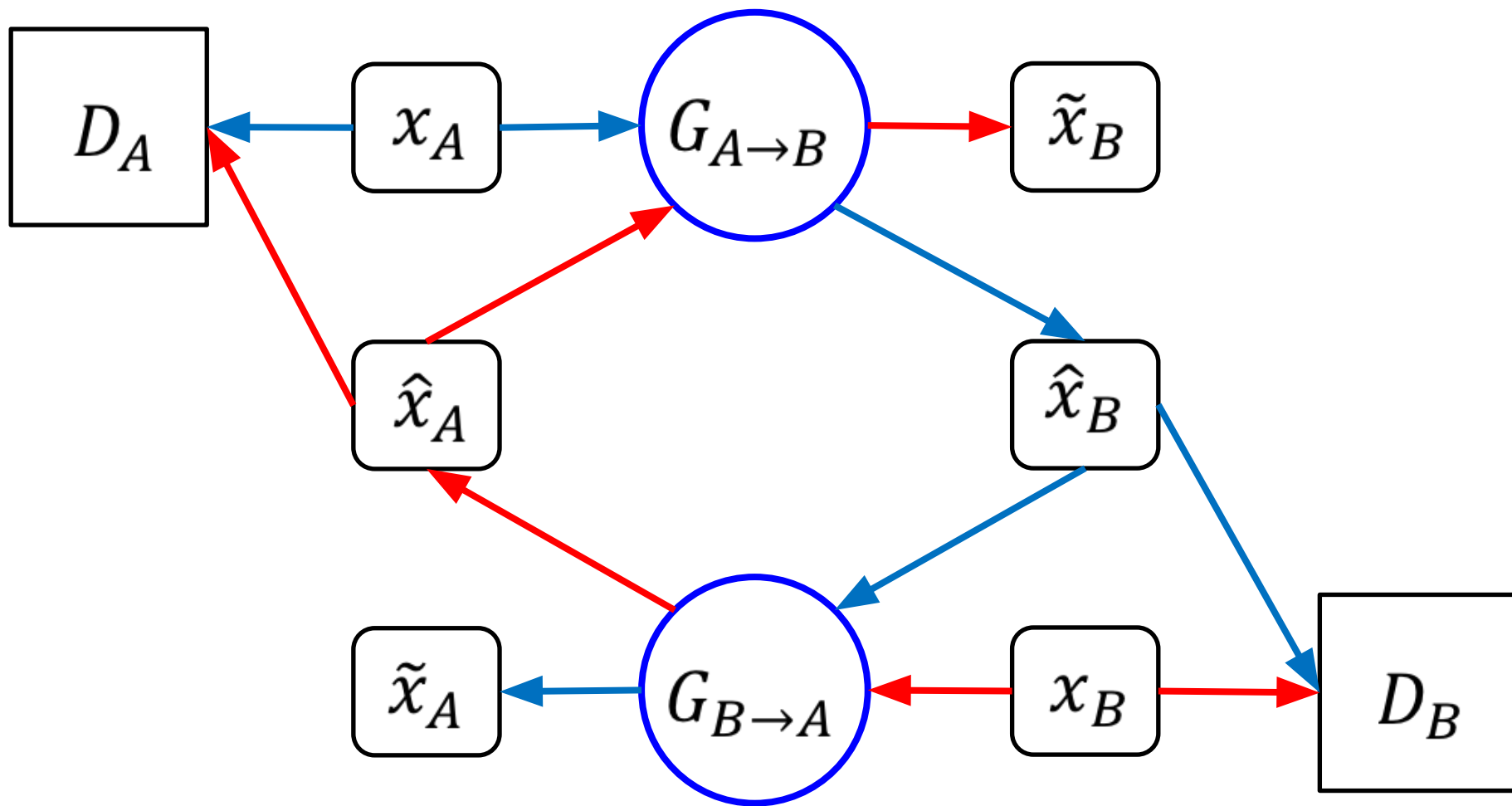
Self-Attention

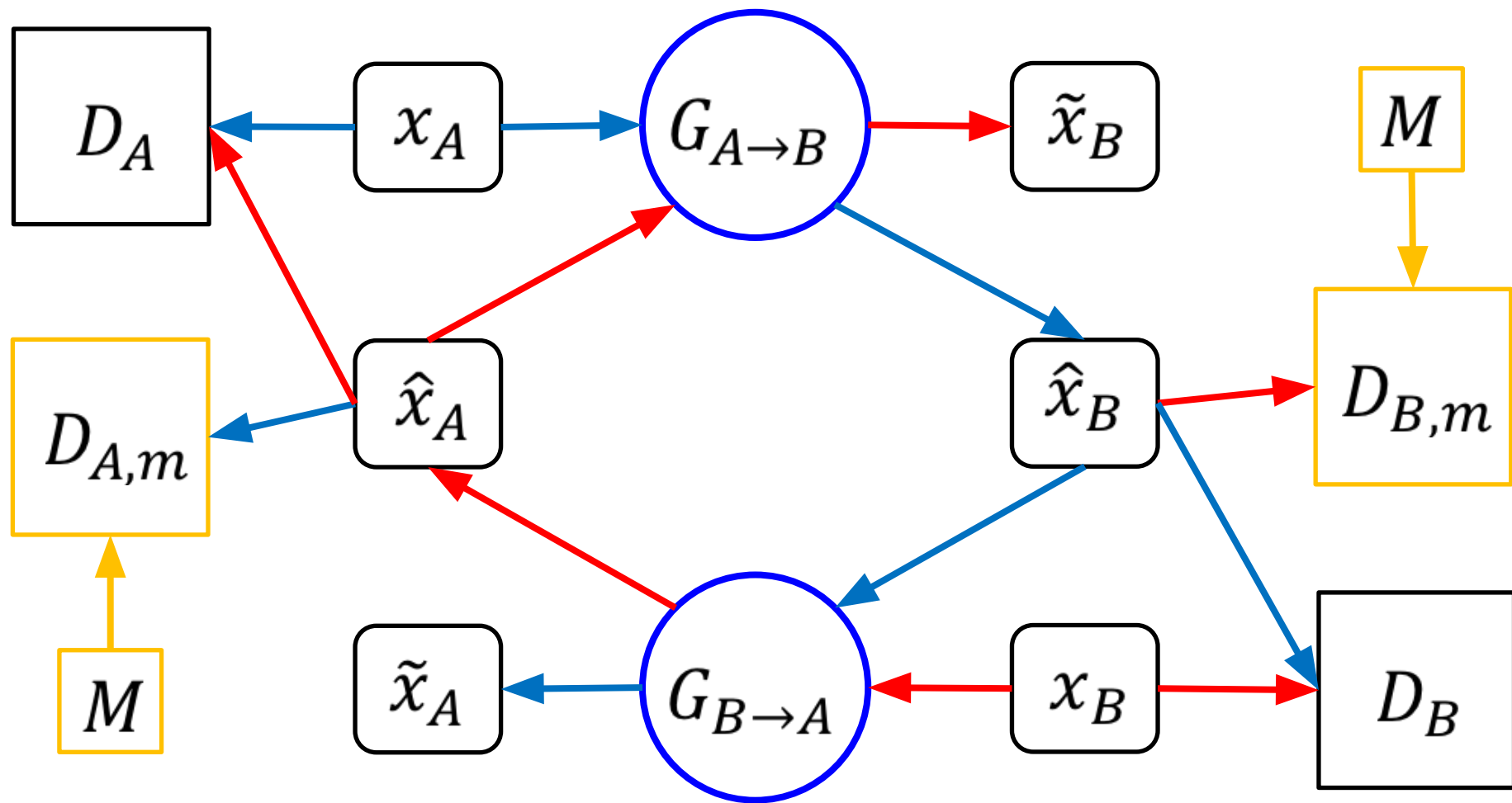
Spectral Normalization

SN + SA

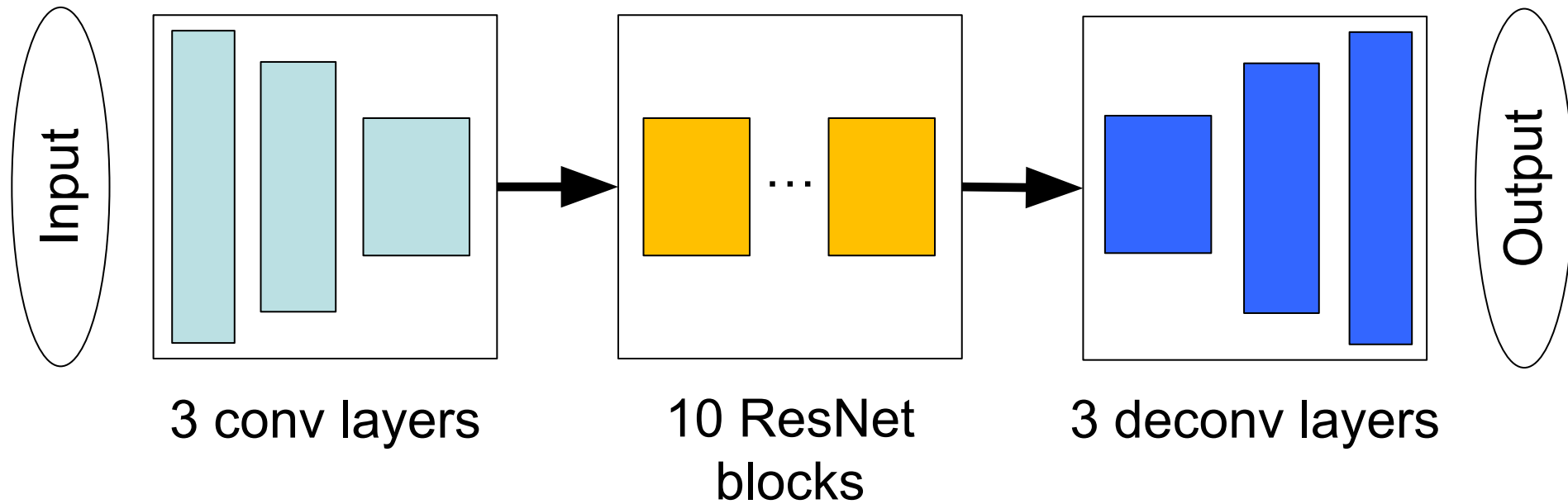


$$s_{hr}^M = \frac{1}{N} \sum_{r=1}^K \#\{\text{rank of } M = r\} \frac{K - r}{K - 1}$$

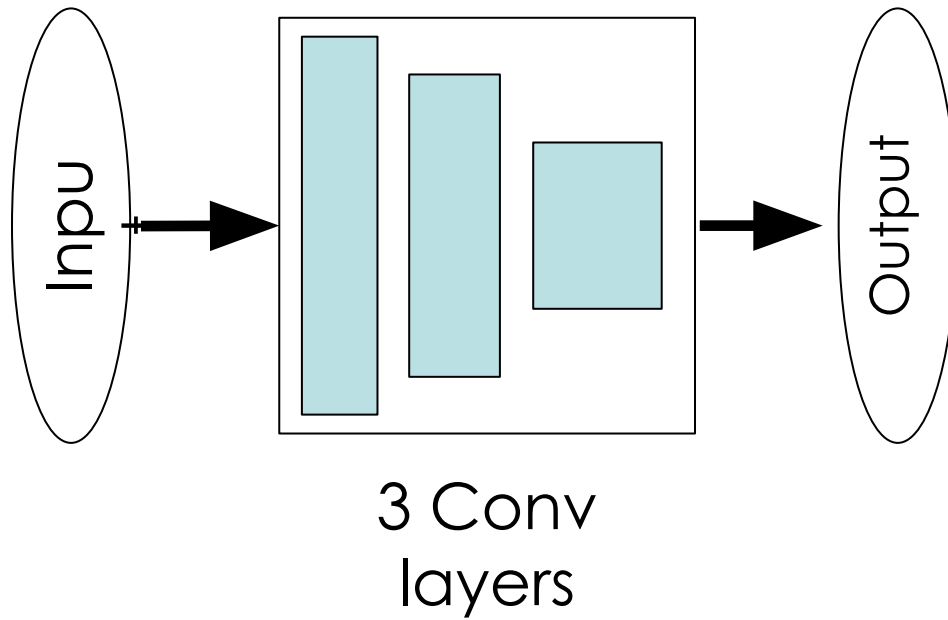




Generators

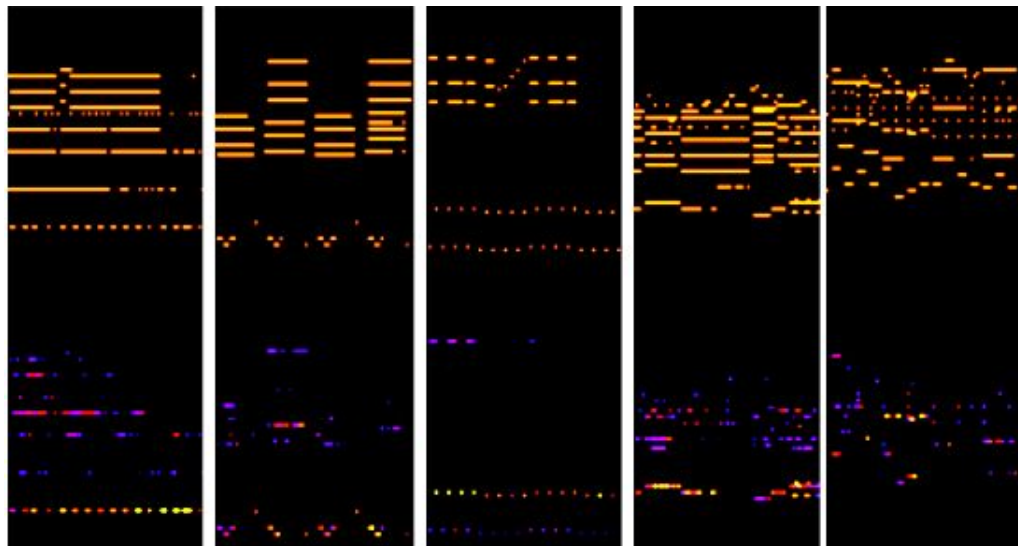


Discriminators

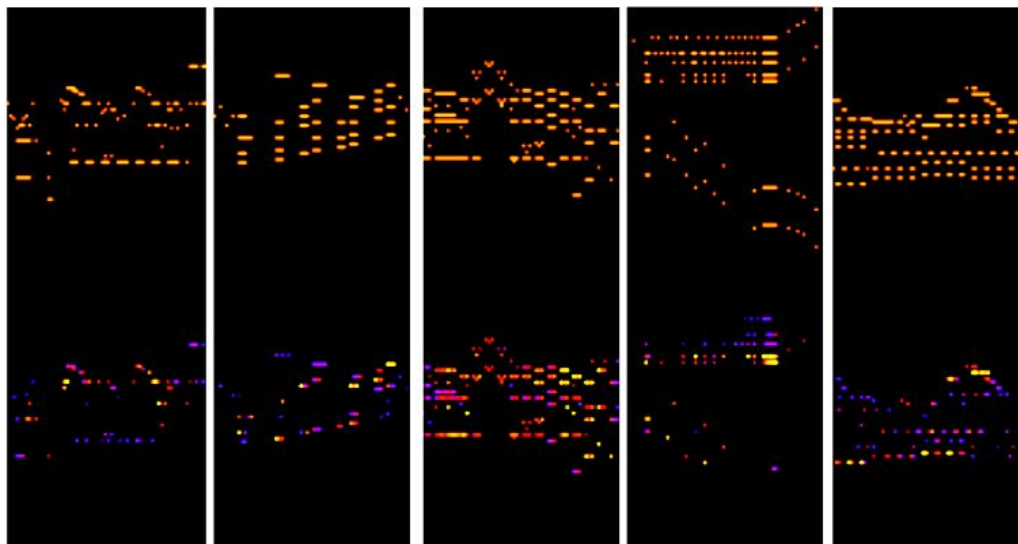


GRADIENT BASED ATTRIBUTION

Jazz samples



Classical samples



Original samples with saliency maps below.