

PAT 498/598 (Winter 2025)

Music & AI

Lecture 14: Language-based Music Generation

Instructor: Hao-Wen Dong



SCHOOL OF MUSIC, THEATRE & DANCE
PERFORMING ARTS TECHNOLOGY
UNIVERSITY OF MICHIGAN

Homework 5: AI Song Contest

- Please listen to the **ten finalists of AI Song Contest 2024**
- **Read the about pages** by clicking the cover arts
- **Answer the following questions** (in 5-10 sentences each)
 - **Which is your favorite song?**
 - Following Q1, **what did they do well?**
 - Following Q1, **what can be improved?**
 - Based on the ten finalists, **what tasks are easy** for current AI in music production?
 - Based on the ten finalists, **what tasks are difficult** for current AI in music production?

Homework 5: AI Song Contest

- Instructions will be released on the [course website](#)
- Please submit your work to [Gradescope](#)
- Due at **11:59pm ET** on **March 14**
- Late submissions: **1 point deducted per day**
- No late submission is allowed a week after the due date

Project

- **Open-ended group project** (group size: 2–3)
 - **Building a new AI music tool** or **Exploring creative & artistic use of AI tools**
- **Milestones**
 - **Pitch:** Mar 19
 - **Presentation:** Apr 21
 - **Final report:** Apr 28
- Due at **11:59pm ET** on the date specified
- **No late submissions!** Submit your work early and update it later.

Project Pitch

- Brief 10-min presentation
 - **Team member introduction**
 - **Topic:** What do you want to work on?
 - **Topic:** Who is the target audience/user/customer/reader?
 - **Methodology:** How are you going to approach it?
 - **Methodology:** What are the tools (programming languages, platforms, plugins, hardware, etc.) that you'll be using?
 - **Expected results:** What are the expected deliverables (e.g., an instrument, a plugin, a web/mobile app, a standalone software, an installation, a performance, a composition)?
 - **Planning:** What are the milestones? What do you expect to achieve by the end of February and March?

Project Pitch

- Send me an email with the following info by **11:59 PM ET on March 19**
 - **Names and U-M IDs of all team members**
 - **Topic:** What do you want to work on?
 - **Topic:** Who is the target audience/user/customer/reader?
 - **Methodology:** How are you going to approach it?
 - **Methodology:** What are the tools (programming languages, platforms, plugins, hardware, etc.) that you'll be using?
 - **Expected results:** What are the expected deliverables (e.g., an instrument, a plugin, a web/mobile app, a standalone software, an installation, a performance, a composition)?
 - **Planning:** What are the milestones? What do you expect to achieve by the end of February and March?

AI Song Contest

AI Song Contest

- Annual international competition showcasing the **creative potential of human-AI co-creativity in the songwriting process**

aisongcontest.com



Yaboi Hanoi – Entering Demons & Gods (2022)



youtu.be/PbrRoR3nEVw

soundcloud.com/yaboi-hanoi/enter-demons-and-gods



Reading: The Making of Entering Demons & Gods (2022)

“It was like a saxophonist trained in classical Thai motifs, who played a special ‘Thai Edition’ saxophone with Phi Nai tunings, had joined the musical conversation. The same was true with the trumpet model and the ขลุ่ย ‘Khlui’ - a flute from Thai, Laos and Cambodian repertoire. I could assemble a **transcultural ensemble** to expand the sonic palette of Thai motifs, whilst adhering to underlying tunings and idiomatic inflections like never before.”

lamtharnhantrakul.github.io/enter-demons-and-gods/



Synthetic Beat Brigade - How would you touch me? (2023)



youtu.be/O4cJ3acEGDw &
drive.google.com/file/d/1QTQ7P3iZI6l0anlwNQ3ewf8g3JjDjesl/view

Synthetic Beat Brigade - How would you touch me? (2023)

- **Ideation:** Spotify API, ChatGPT, Facebook Llama, Google Bison
- **Lyrics:** ChatGPT 2, Genius API
- **Composition:** AI Drummachine, Mofi, Tonetransfer, This patch does not exist, Albeatz, BaiscPitch, Magenta, AIVA, MuseNet
- **Vocals:** Soundly Voice Designer, Vocal Remove, Voice characteristics
- **Mastering:** Landr
- **Cover art, bandart:** Midjourney
- **Clip:** ComfyUI for Stable Diffusion + ControlNet

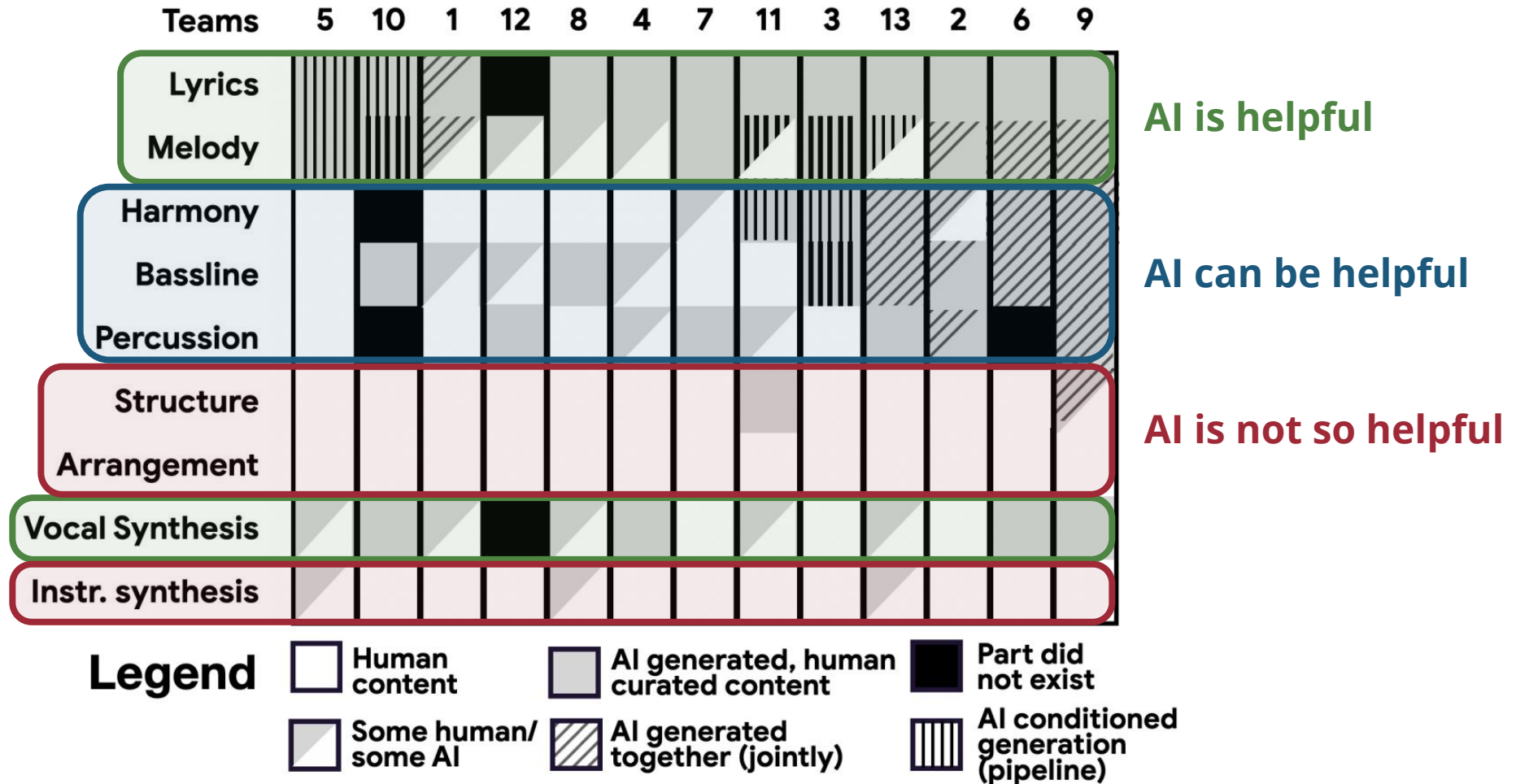
Reading: The Making of How would you touch me? (2023)

“This project is a collaboration between **Artificial Intelligence (AI) enthusiasts in four fields: artist management, music and post-production, tech, and creative**. In contrast, the majority of the music industry sees AI as a threat. Our team understands that these technological advances will have a significant impact on how we produce music. Because of this, we have decided to **use AI for every step of the production process**. From ideation to creating the lyrics to producing the music.”

drive.google.com/file/d/1QTQ7P3iZI6l0anlwNQ3ewf8g3JjDjesl/view

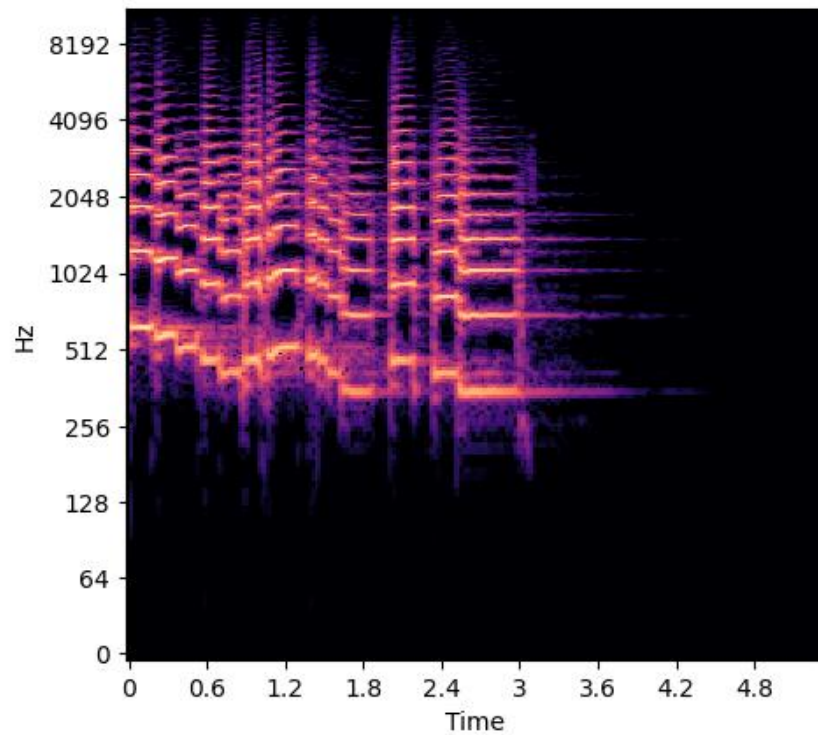


How can AI Augment Human Creativity?

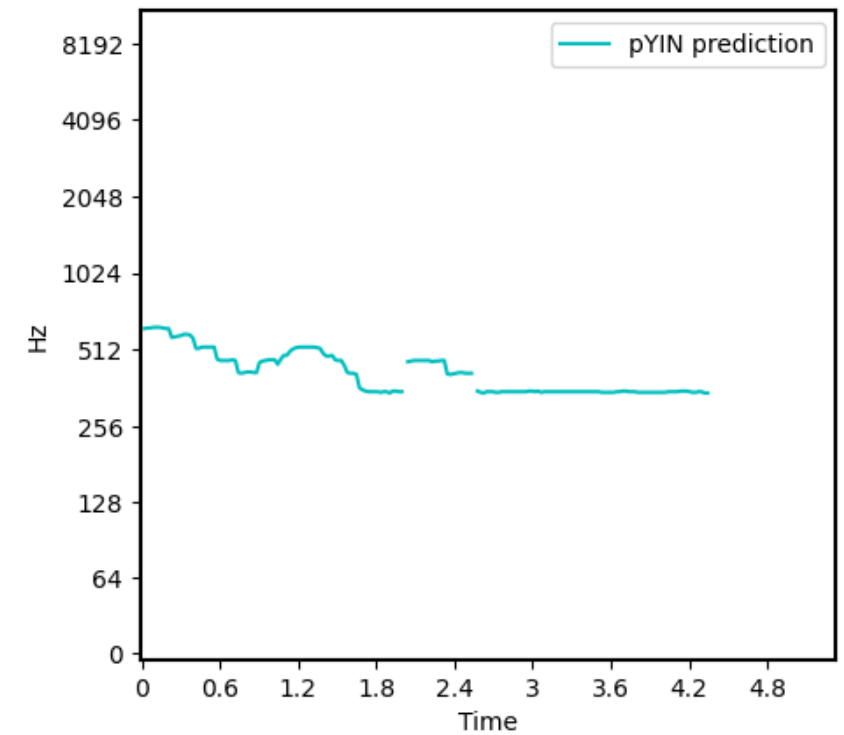


(Source: Huang et al., 2020)

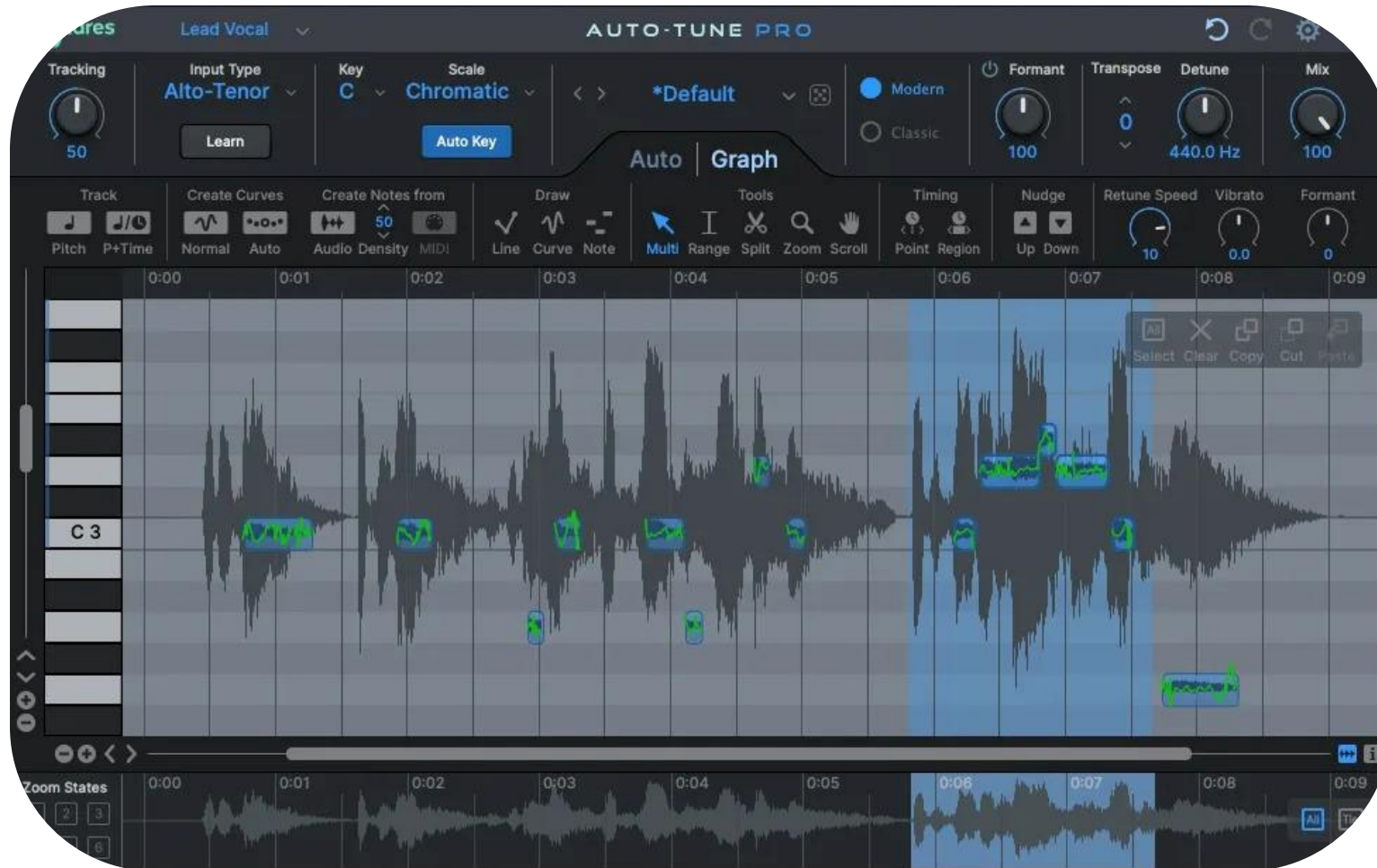
(Recap) Fundamental Frequency (F0) Estimation



→ **F0 Estimation** →

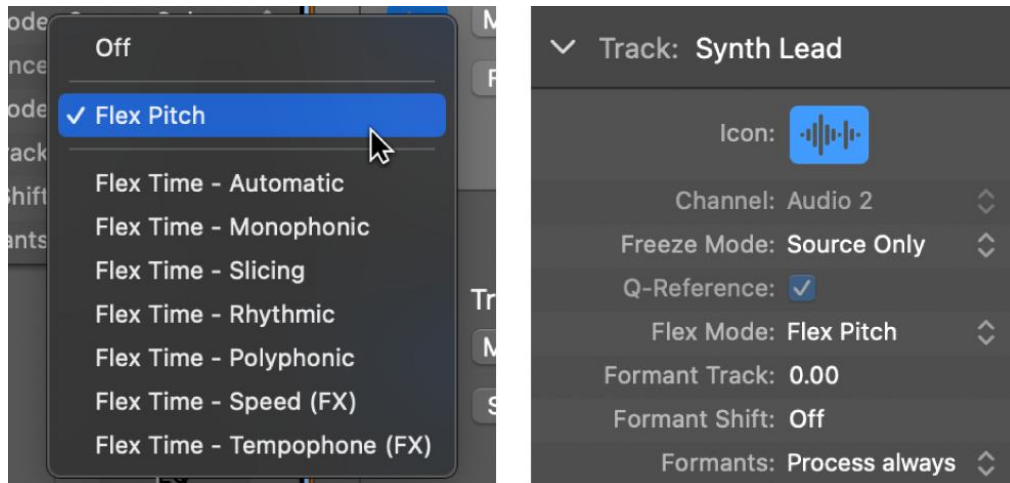


(Recap) Auto-tune Pro

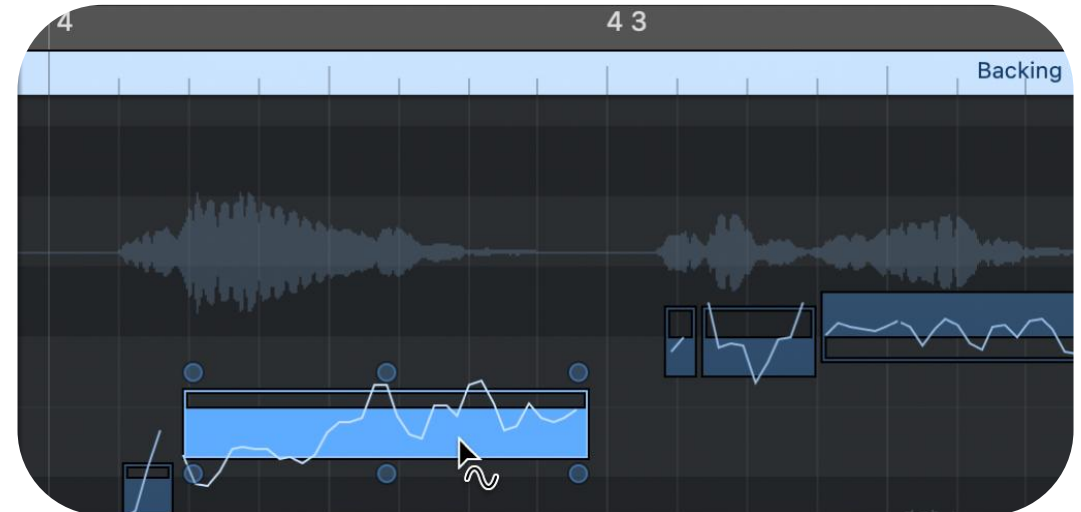


(Source: Antares Audio Technologies)

(Recap) Pitch Correction in Logic Pro

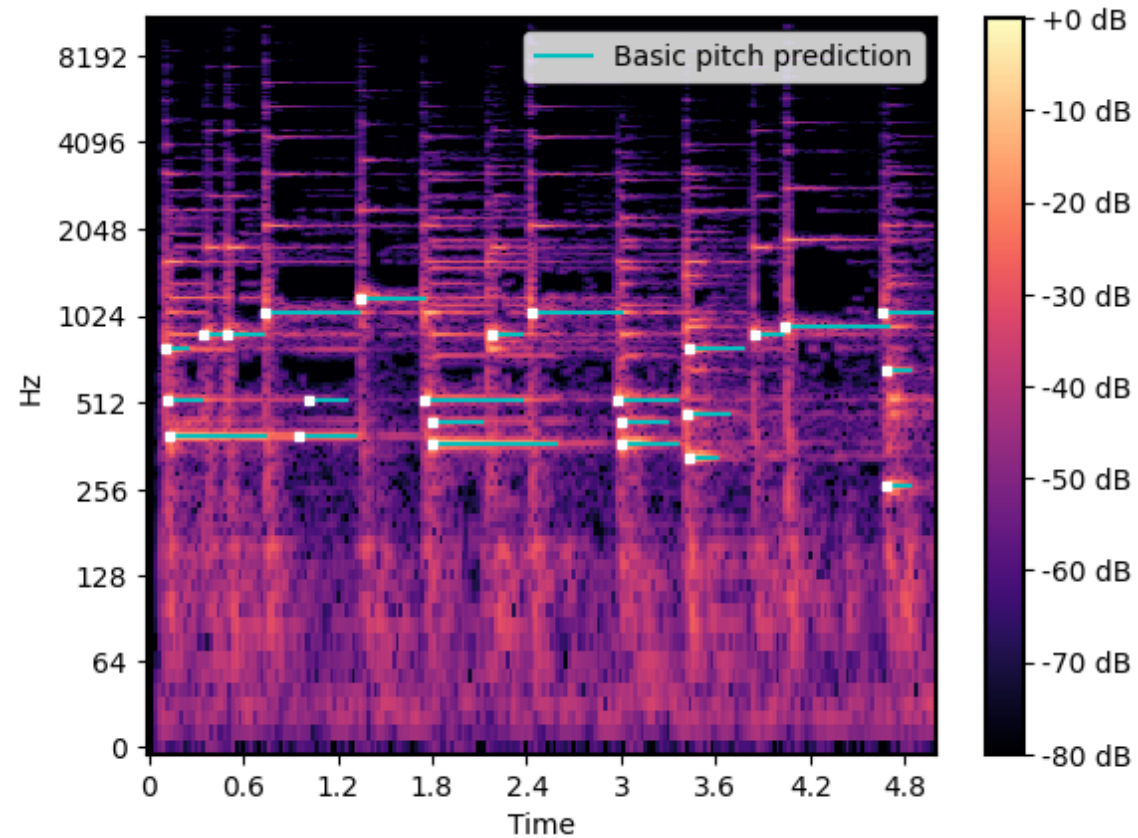


(Source: Logic Pro User Guide)



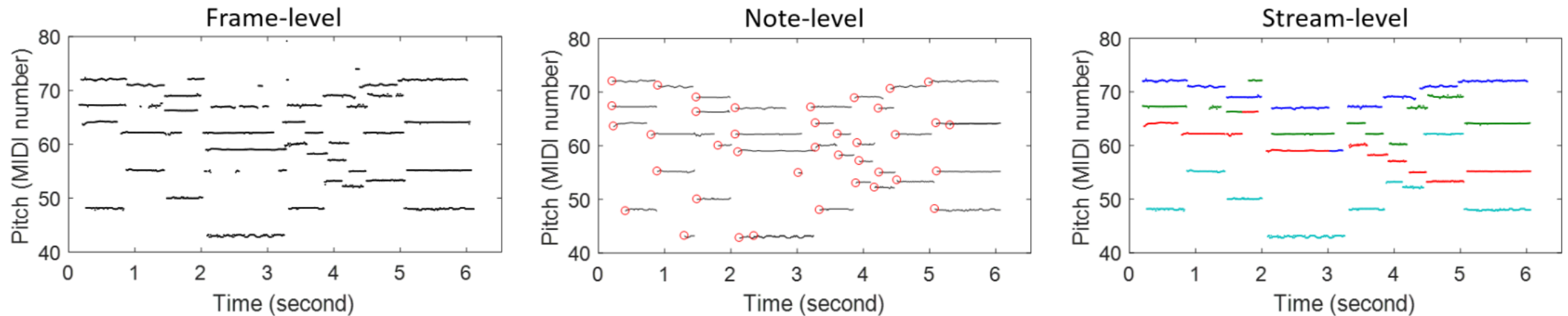
(Source: Logic Pro User Guide)

(Recap) Polyphonic F0 Estimation



basicpitch.spotify.com

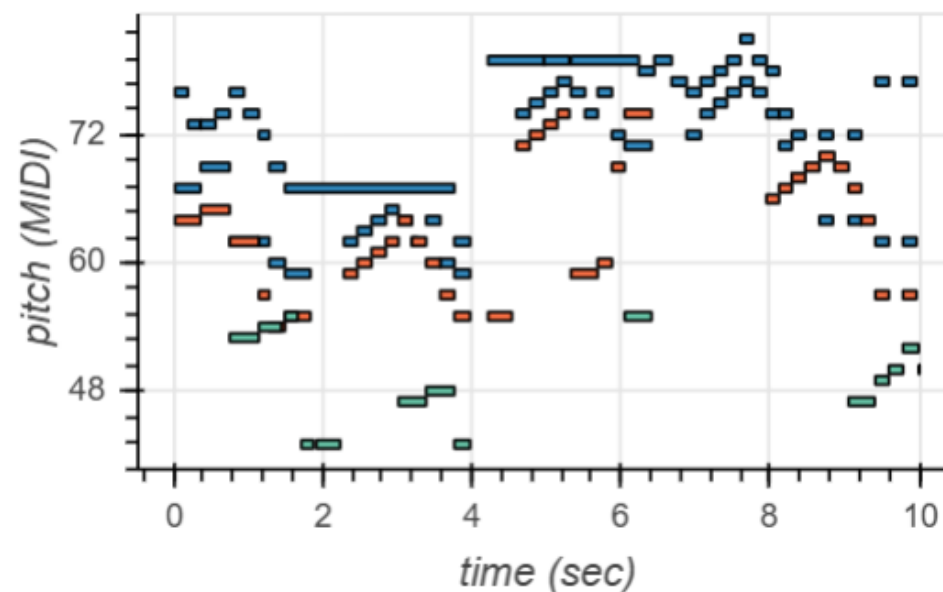
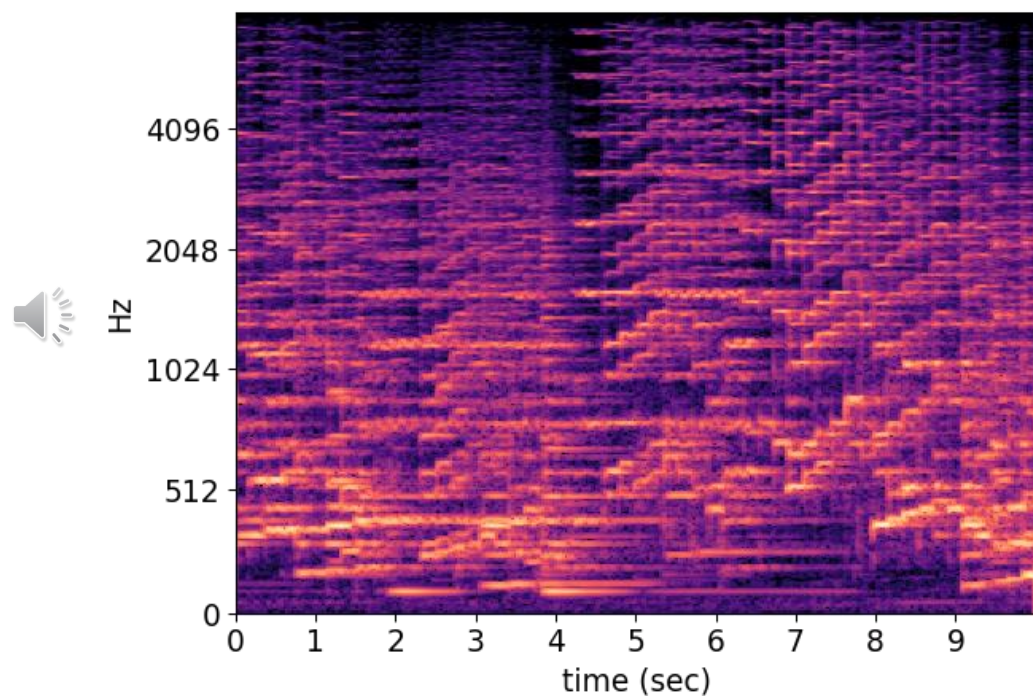
(Recap) F0 Estimation vs Music Transcription



(Source: Benetos et al., 2019)

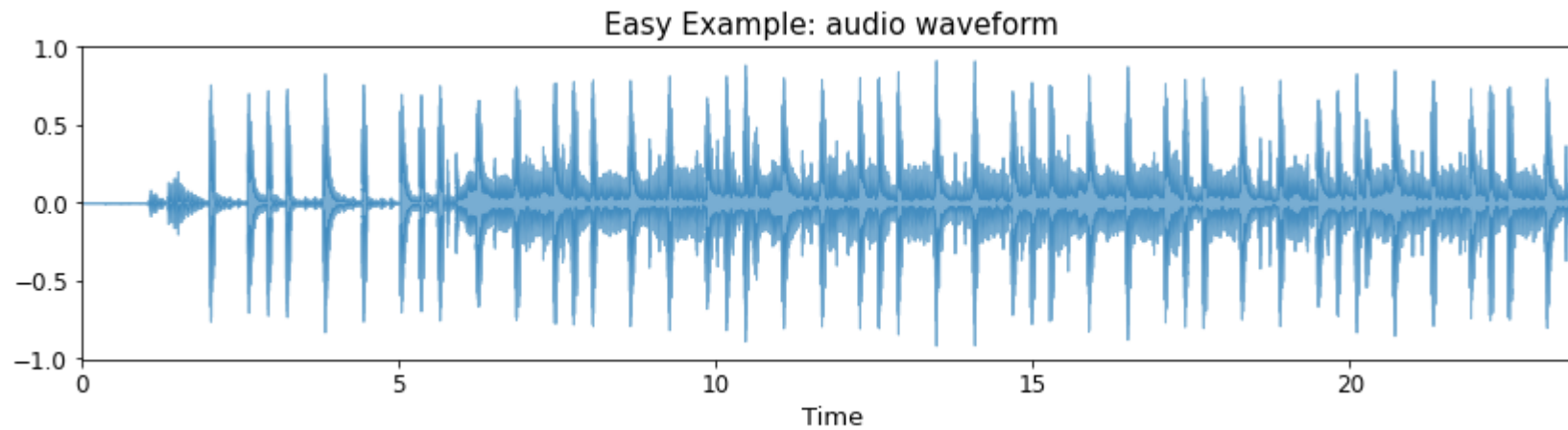
(Recap) Multitrack Transcription Models

- **MT3** (Gardner et al., 2022)
 - github.com/magenta/mt3



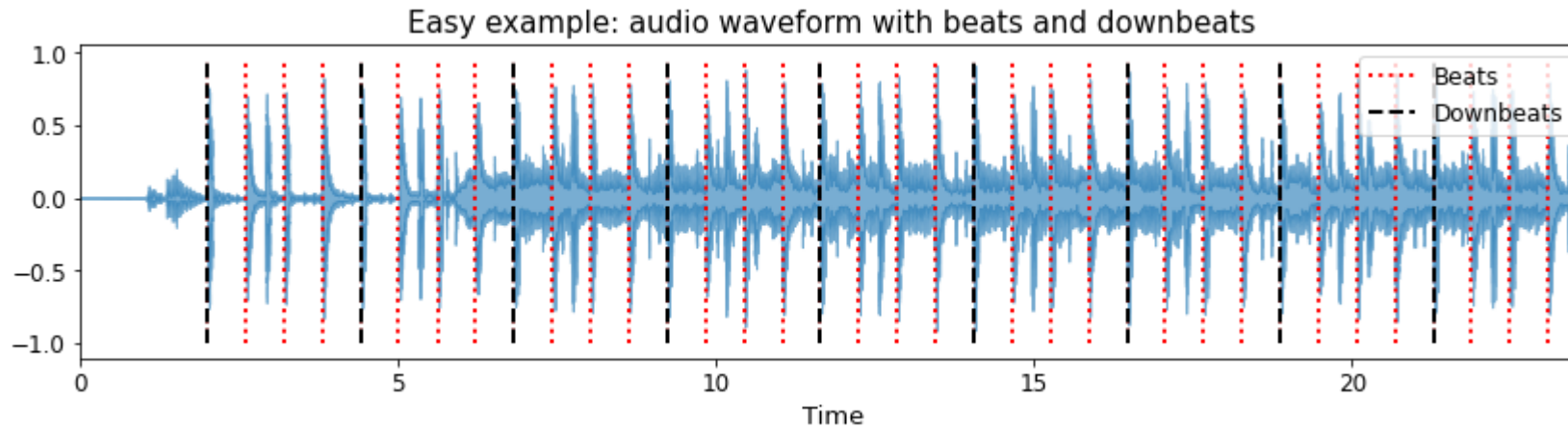
(Source: Gardner et al., 2022)

(Recap) Beat & Downbeat Estimation



(Source: Davies et al., 2021)

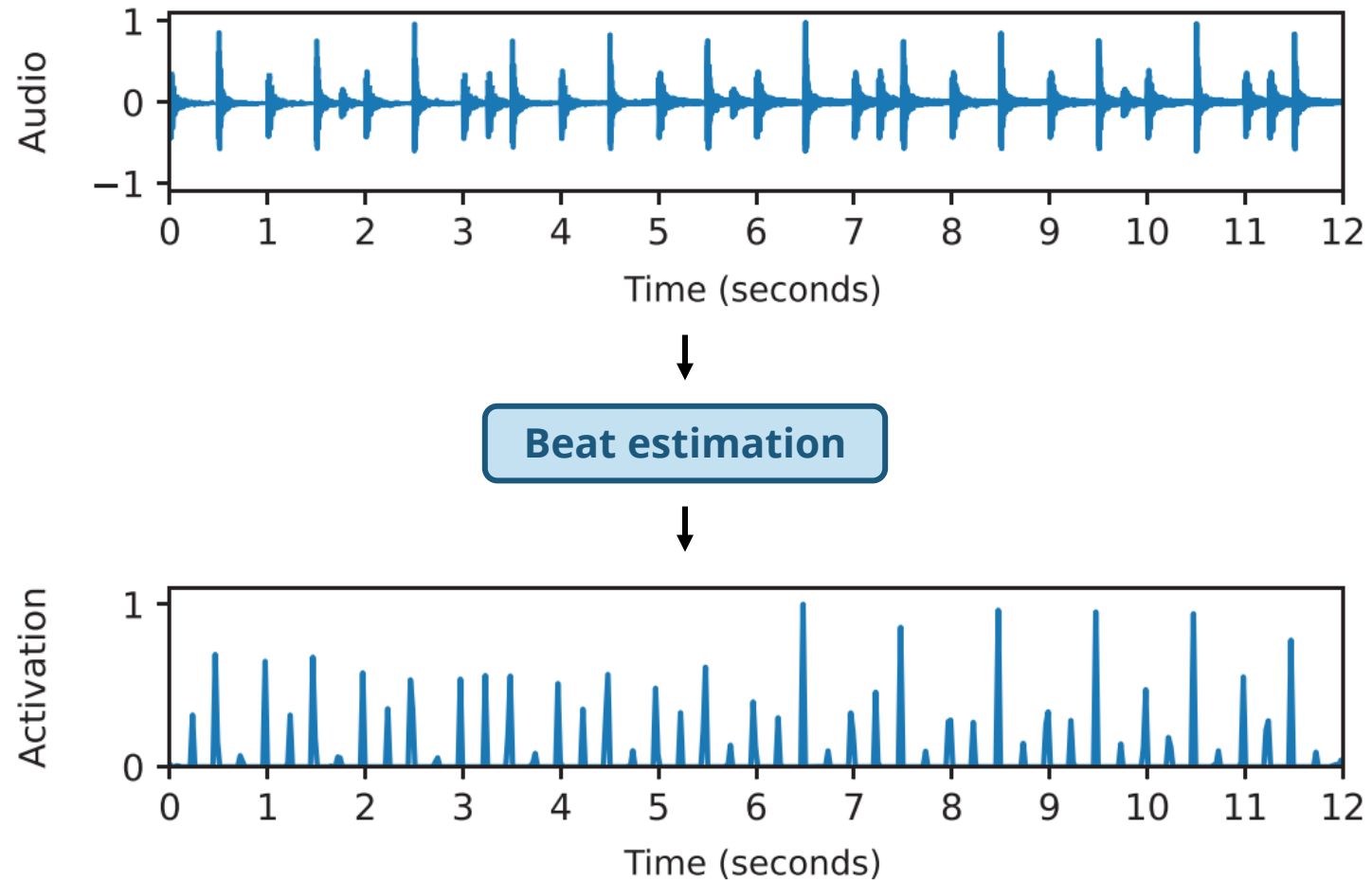
(Recap) Beat & Downbeat Estimation



tatum
beat
downbeat

(Source: Davies et al., 2021)

Beat & Downbeat Estimation



(Source: Meier et al., 2024)

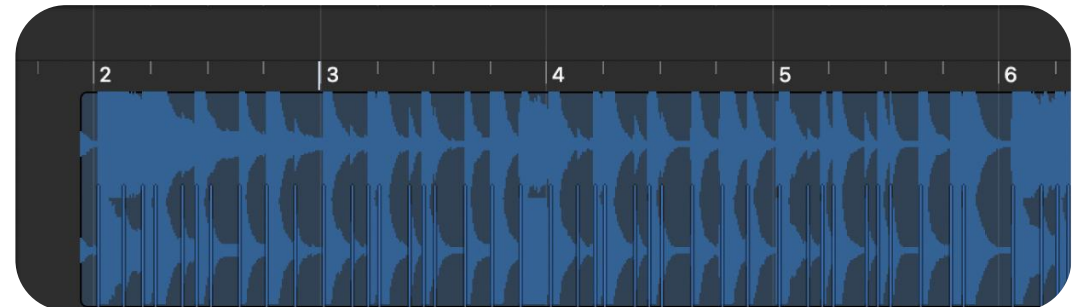
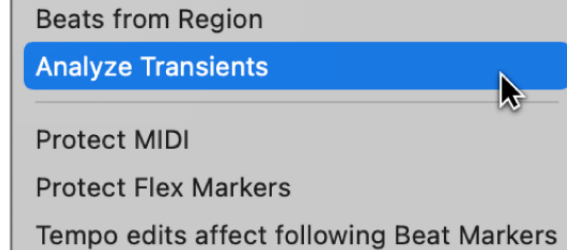
(Recap) Beat Tracking in Pro Tools & Logic Pro

Beat Detective in Pro Tools



(Source: Logic Pro User Guide)

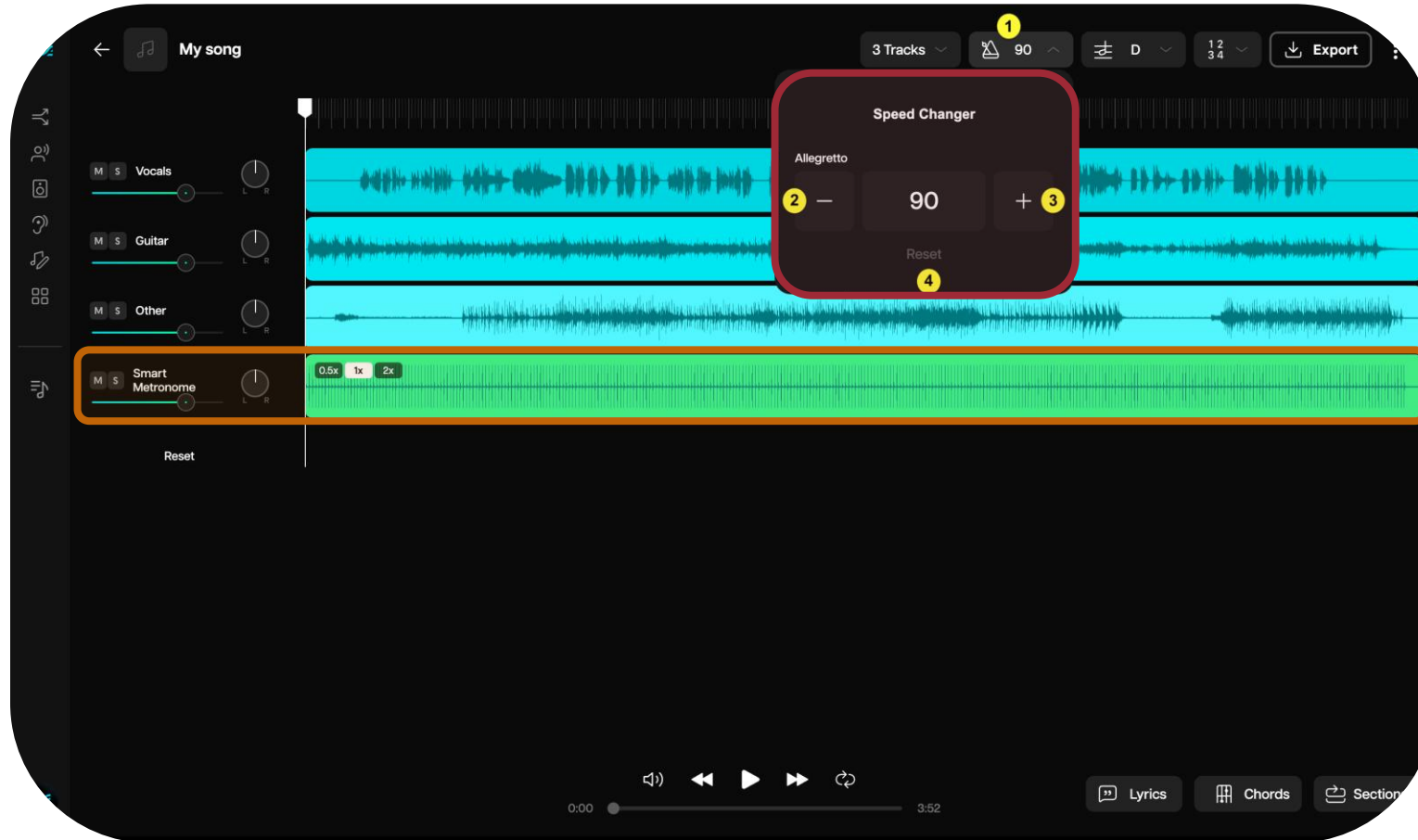
Beat Mapping in Logic Pro



(Source: Logic Pro User Guide)

(Recap) Tempo Estimation & Beat Tracking in Moises

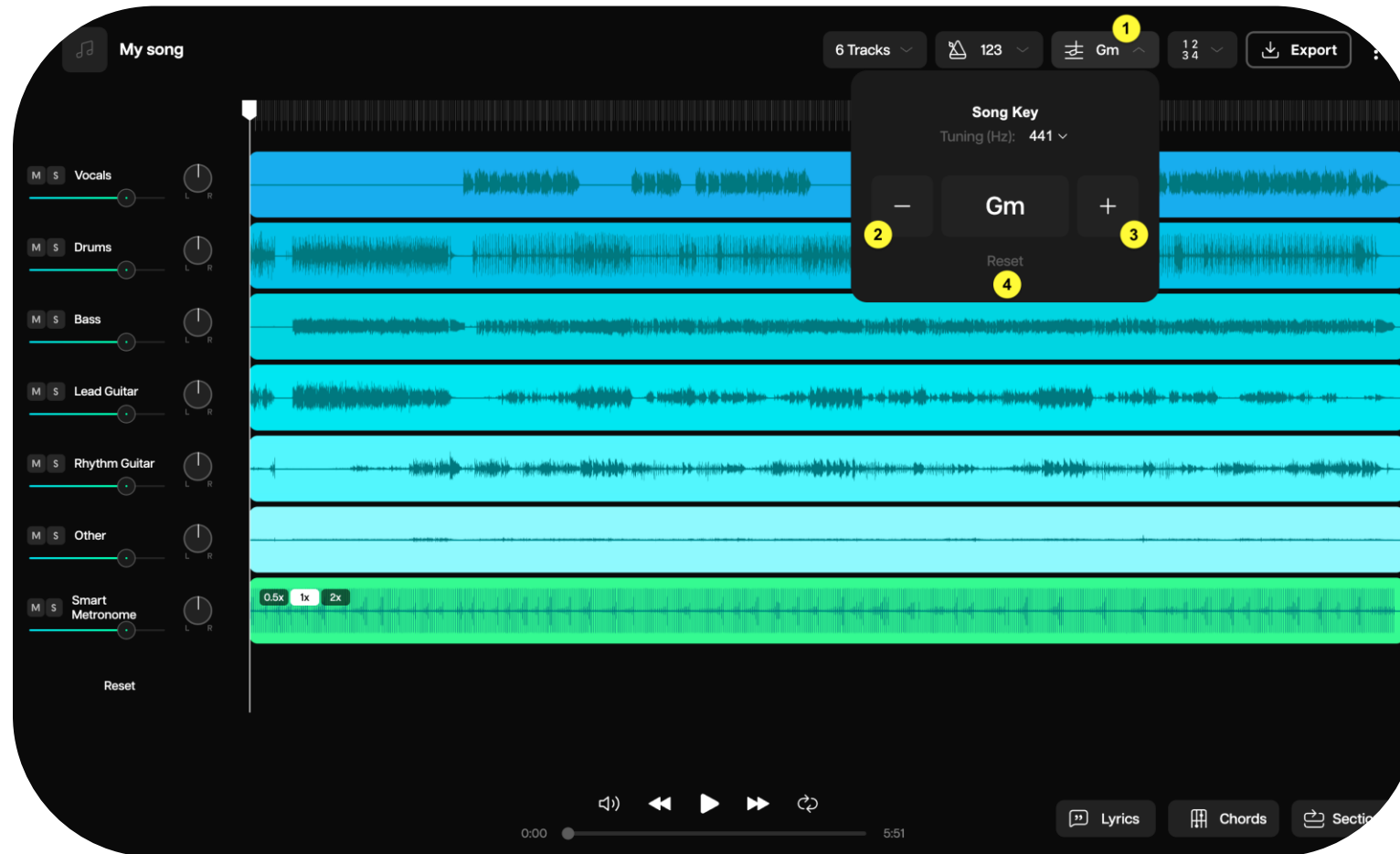
Tempo estimation



Beat tracking

(Source: Moises)

(Recap) Key Detection in Moises



(Source: Moises)

(Recap) Structure Analysis

Music segmentation

The musical score is divided into segments labeled A1, A2, B1, B2, C, A3, B3, B4, and D. The score includes dynamic markings such as *f*, *ff*, *p*, *ff marc.*, *sf*, *p poco rit.*, *in tempo*, *p poco ritard.*, *p legg.*, *p in tempo*, *poco ritard.*, *in tempo*, *p*, *ff marcato*, *p poco ritard.*, and *ff*. The tempo markings include *Allegro.* and *Vivace.*

Figure 4.5 following [Müller, FMP, Springer 2015]

(Source: Müller & Zalkow, 2019)

Hierarchical music segmentation

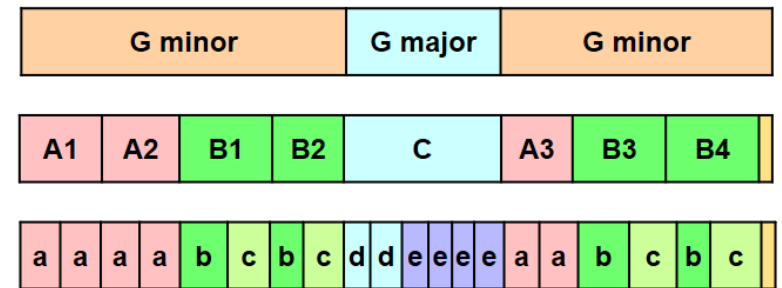
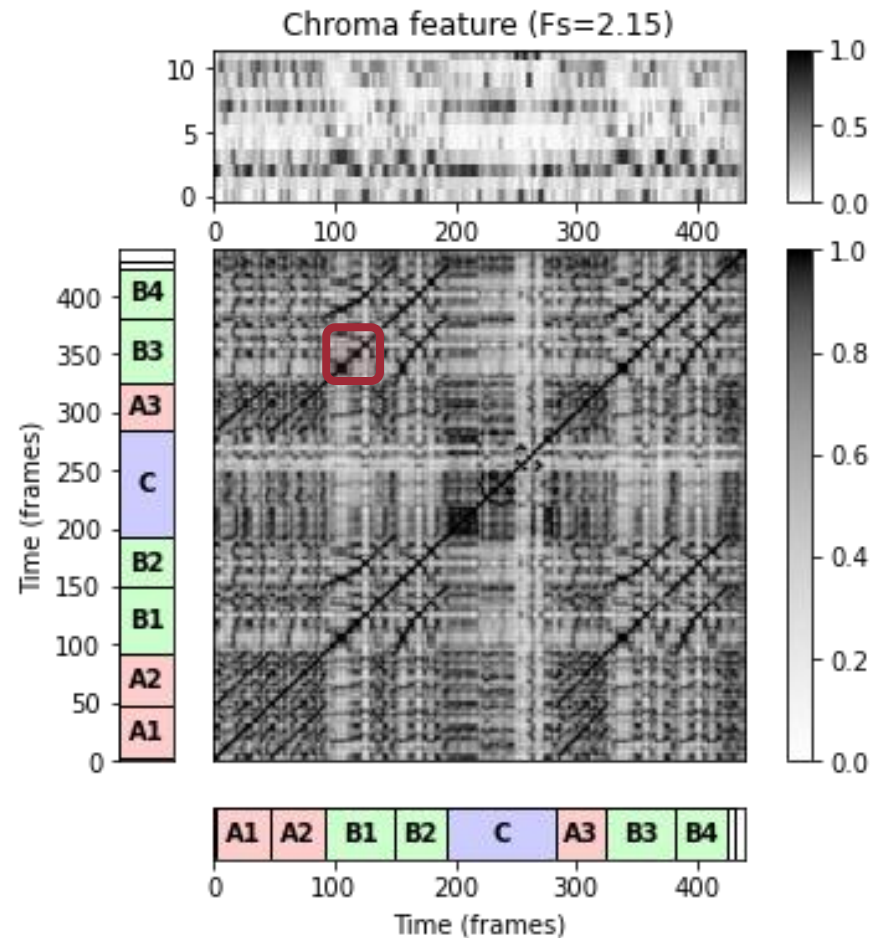


Figure 4.28 from [Müller, FMP, Springer 2015]

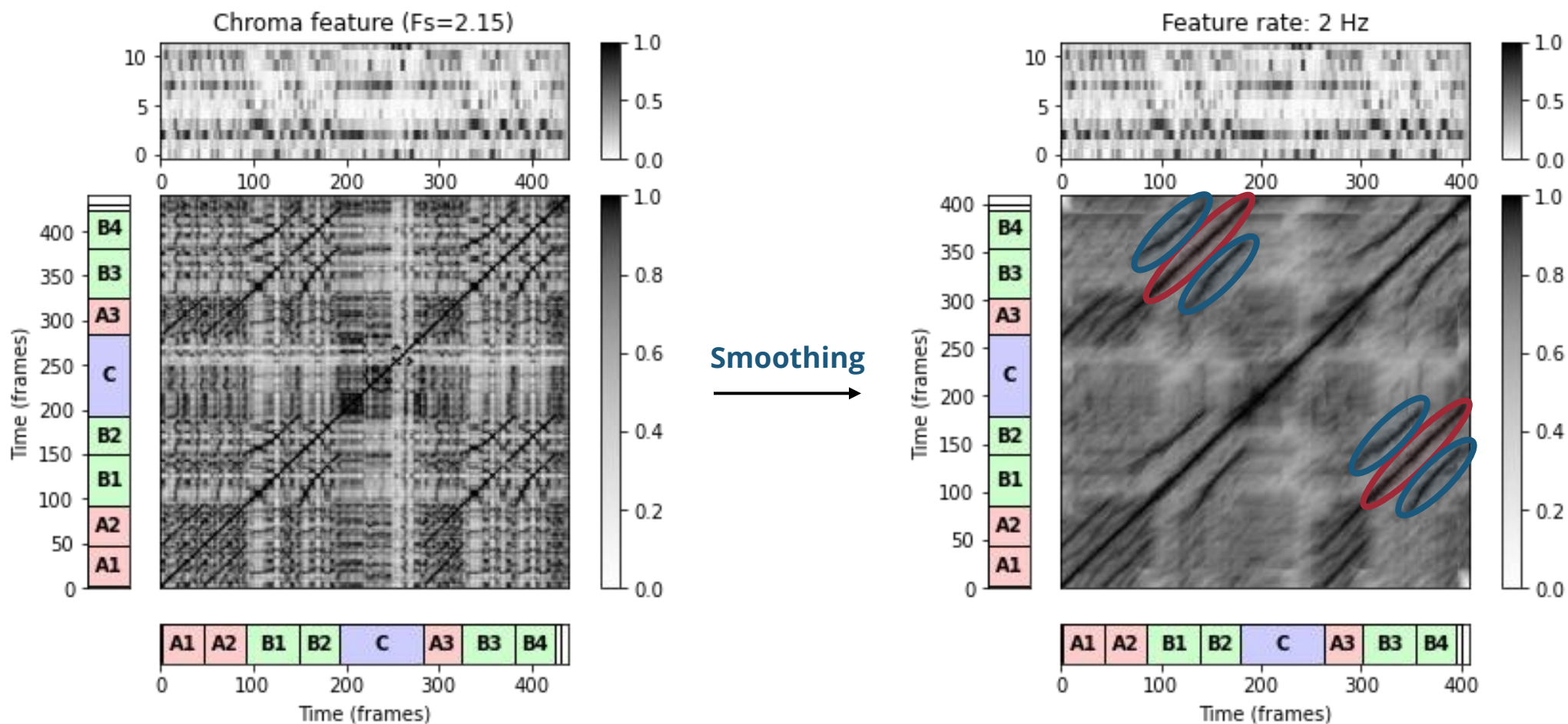
(Source: Müller & Zalkow, 2019)

(Recap) Self-Similarity Matrices (SSMs)



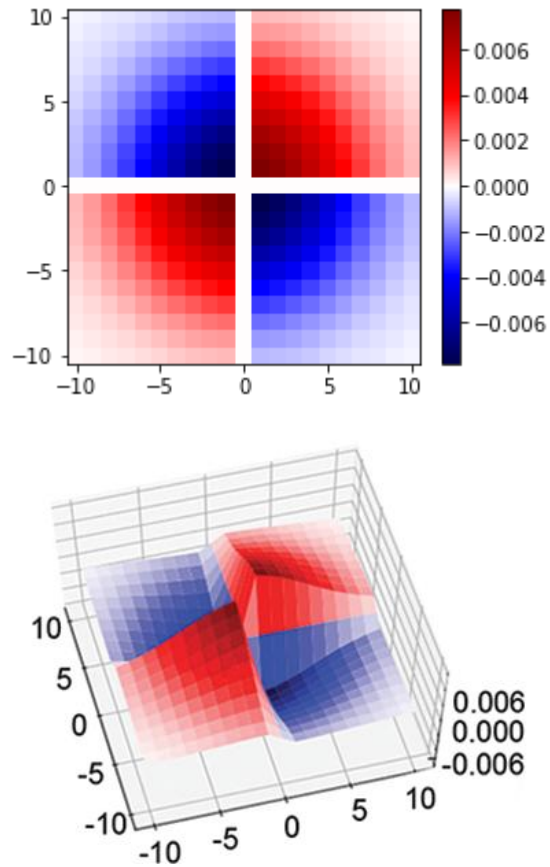
(Source: Müller & Zalkow, 2019)

(Recap) Self-Similarity Matrices (SSMs)



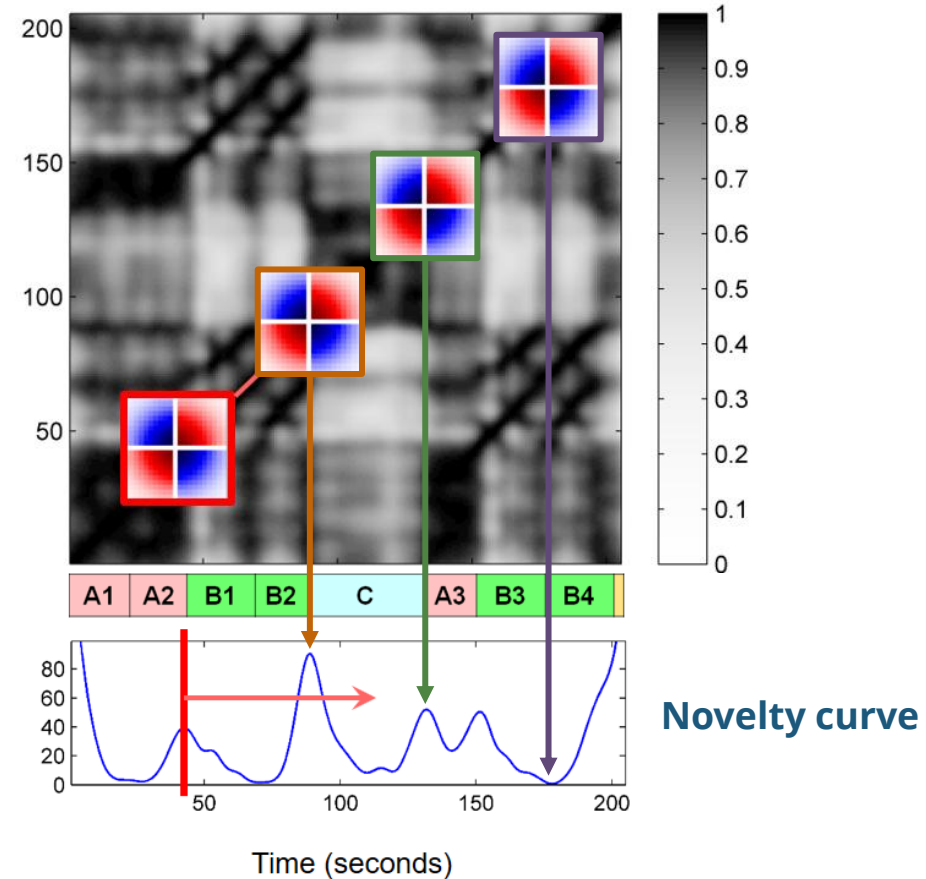
(Source: Müller & Zalkow, 2019)

(Recap) Self-Similarity Matrices (SSMs)



(Source: Müller & Chiu, 2024)

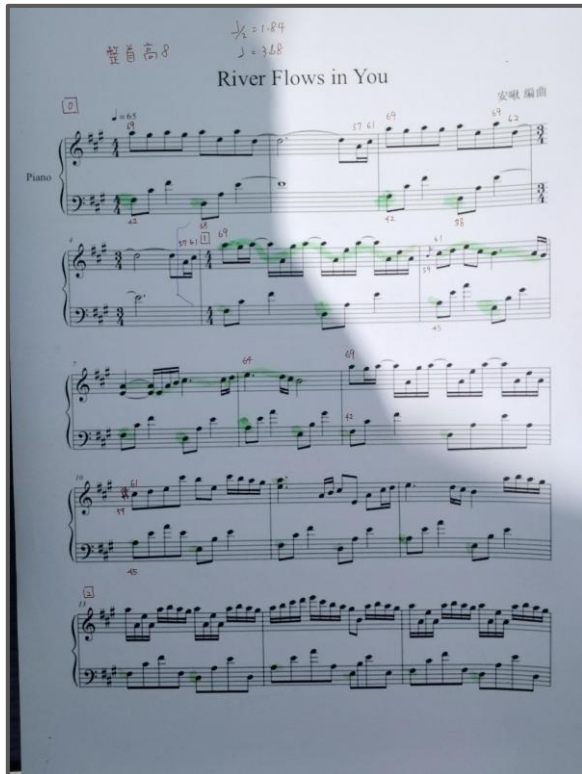
Figure 4.24 from [Müller, FMP, Springer 2015]



(Source: Müller & Zalkow, 2019)

(Recap) Optical Music Recognition (OMR)

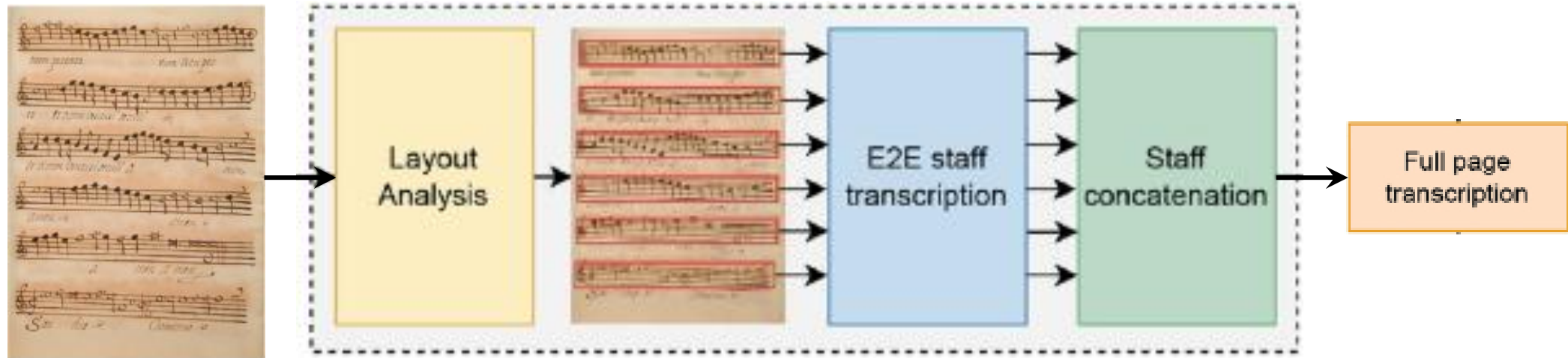
- **Goal:** Convert **scanned sheet music** into **digital musical notation**



Optical Music Recognition

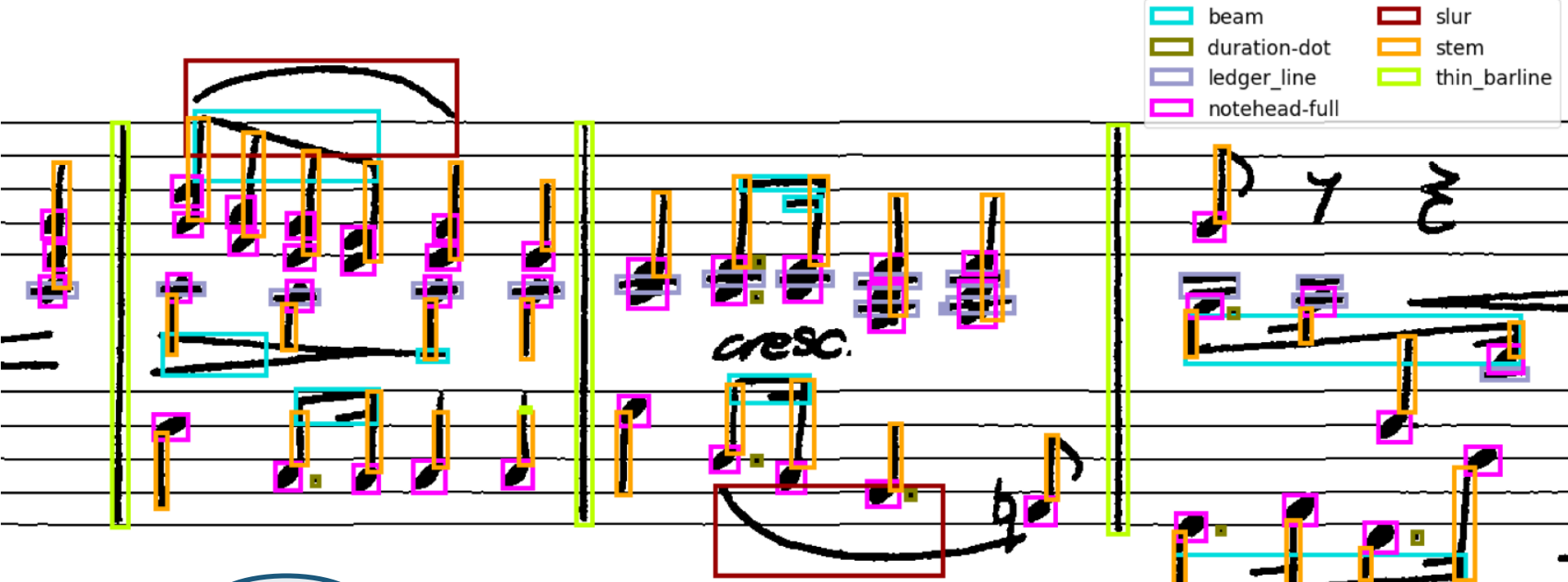


(Recap) Common Pipeline of OMR Systems



(Source: Calvo-Zaragoza et al., 2018)

(Recap) Musical Object Recognition



A musical score on three staves with various annotations. A legend in the top right corner identifies the following objects: beam (cyan), duration-dot (olive), ledger_line (purple), notehead-full (magenta), slur (red), stem (orange), and thin_barline (yellow-green). The score includes a 'cresc.' marking and a 'b' (bass clef) marking. Two blue circles with arrows point to a vertical barline and a slur, with the text 'From the staff below!' below them.

beam	slur
duration-dot	stem
ledger_line	thin_barline
notehead-full	

From the staff below!

(Source: Pacha et al., 2018)

(Recap) Open-source OMR Software: Oemer

Composed by Tomohiko Kira
Arranged by Asimenez
Transcribed by maDwaZz

mp
mf
f
cantabile

Tabi

Transcribed by Oemer

1
transcribed by Oemer

3
transcribed by Oemer

github.com/BreezeWhite/oemer

Symbolic Music Generation

Four Paradigms of Music Generation



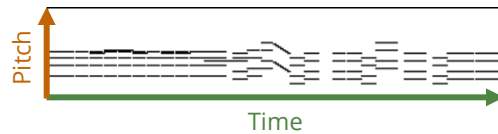
Symbolic music generation

Text-based

Image-based

```
Program_change_0,  
Note_on_60, Time_shift_2, Note_off_60,  
Note_on_60, Time_shift_2, Note_off_60,  
Note_on_76, Time_shift_2, Note_off_67,  
Note_on_67, Time_shift_2, Note_off_67,  
...
```

MIDI



Piano roll



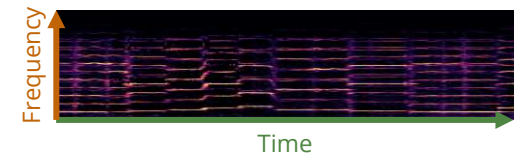
Audio-domain music generation

Time series-based

Image-based



Waveform



Spectrogram

Today, we also have many **latent-space based systems!**

Topics of Symbolic Music Generation

Unconditional

Symbolic music generation

- $\emptyset \rightarrow$ melody
- $\emptyset \rightarrow$ lead sheet
- $\emptyset \rightarrow$ sheet music

Melody
& chords

Today's topic!

Conditional

Automatic arrangement

- Melody \rightarrow lead sheet
- Melody \rightarrow multitrack
- Lead sheet \rightarrow multitrack
- Solo \rightarrow multitrack
- Multitrack \rightarrow simple version

Performance rendering

- Sheet music \rightarrow performance

Improvisation systems

- Performance \rightarrow performance

Multimodal

X-to-music generation

- Text \rightarrow sheet music
- Video \rightarrow sheet music
- X \rightarrow sheet music

Two Paradigms of Symbolic Music Generation



Text-based

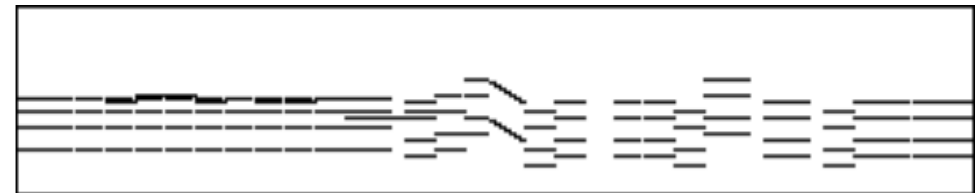
- Treat music like **text**
- Sharing models with **natural language processing (NLP)**
 - RNNs, LSTMs, Transformers, etc.

Today's topic!

```
Program_change_0,  
Note_on_60, Time_shift_2, Note_off_60,  
Note_on_60, Time_shift_2, Note_off_60,  
Note_on_76, Time_shift_2, Note_off_67,  
Note_on_67, Time_shift_2, Note_off_67, ...
```

Image-based

- Treat music like **images**
- Sharing models with **computer vision (CV)**
 - GANs, VAEs, diffusion models, etc.



Generating Music like Languages

Large Language Models (LLMs)

- The models behind ChatGPT!

SA

You

What's so cool about **AI for music**? Give me a brief answer



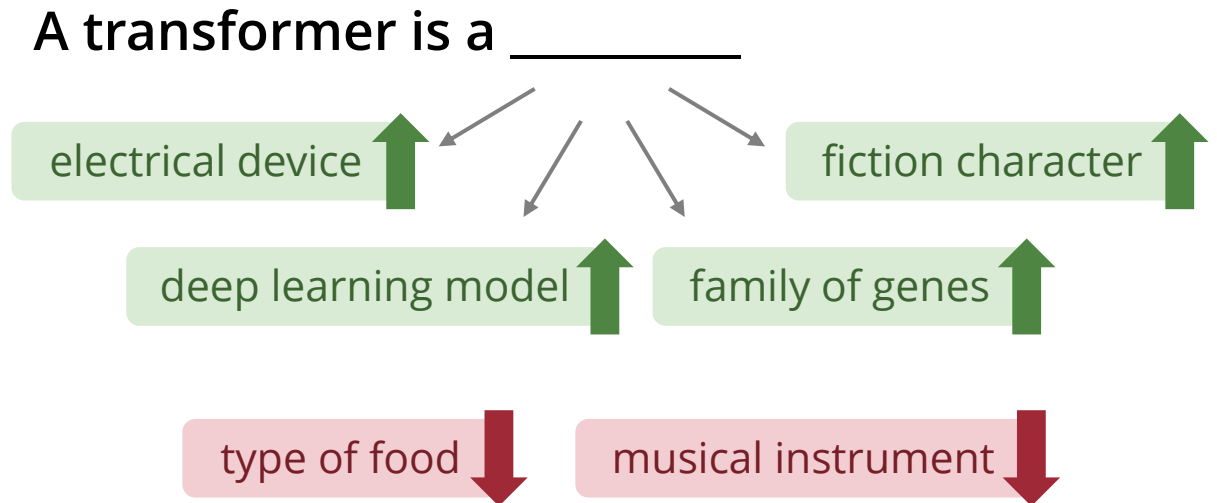
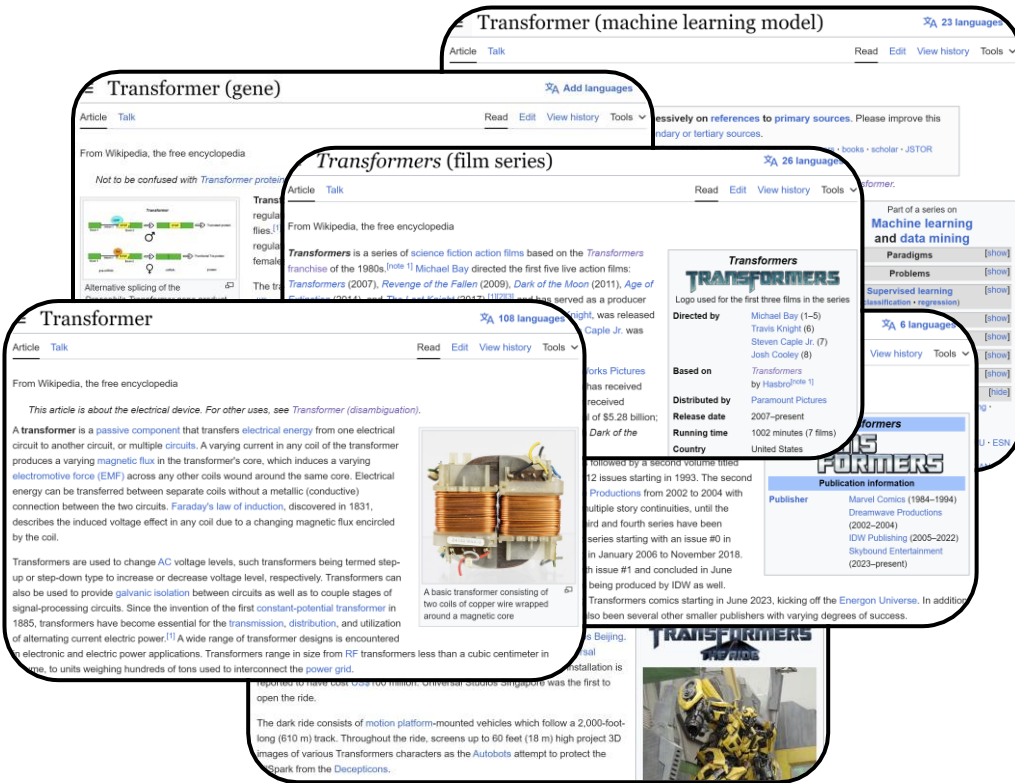
ChatGPT

Word-by-word generation

AI in music is cool because it can compose original pieces, provide personalized recommendations, automate music production tasks, enhance creativity for artists, enable interactive performances, analyze music trends, and even create virtual artists or bands, expanding the possibilities in music creation and enjoyment.

Language Models

- Predicting the next word given the past sequence of words



Language Models (Mathematically)

- A class of machine learning models that **learn** the next word probability

$$P(x_i \mid x_1, x_2, \dots, x_{i-1})$$

Next word Previous words

$P(\text{electrical} \mid \text{A transformer is a})$	↑
$P(\text{character} \mid \text{A transformer is a})$	↑
$P(\text{gene} \mid \text{A transformer is a})$	↑
$P(\text{model} \mid \text{A transformer is a})$	↑
$P(\text{food} \mid \text{A transformer is a})$	↓
$P(\text{musical} \mid \text{A transformer is a})$	↓

Language Models – Generation

- How do we generate a new sentence using a trained language model?

A transformer is a

→ Model → deep

A transformer is a deep

→ Model → learning

A transformer is a deep learning

→ Model → model

A transformer is a deep learning model

→ Model → introduced

A transformer is a deep learning model introduced

→ Model → in

A transformer is a deep learning model introduced in

→ Model → 2017

Designing a Machine-readable Music Language

- How can we “represent” music in a way that machines understand?
 - Musical representation is a key component of a music generation system
- Why not using sheet music “images” directly?
 - Machines still have a hard time reading sheet music
 - A challenging task known as “optical music recognition” (OMR)
- Examples:
 - ABC notation
 - MIDI



ABC Notation-based Representation

(Recap) What is this song in ABC notation?

```
CCGG | AAG2 | FFEE | DDC2 : |  
| : GGFF | EED2 | GGFF | EED2 |  
CCGG | AAG2 | FFEE | DDC2 : |
```

Twinkle, twinkle, little star!

(Recap) An Example of ABC Notation

Ah! vous dirai-je, maman
(Twinkle, twinkle, little star)

anon. (France)

♩ = 120

Metadata

```
X:571
T:Ah! vous dirai-je, maman
T:(Twinkle, twinkle, little star)
C:anon.
O:France
R:Nursery song
M:C Meter
L:1/4 Unit note length (temporal resolution)
Q:120 Tempo
K:C Key
CCGG|AAG2|FFEE|DDC2:|
|:GGFF|EED2|GGFF|EED2|
CCGG|AAG2|FFEE|DDC2:|
```

Example System: Folk RNN (Sturm et al., 2015)

- **Data**

- Collections of folk tunes

- **Representation**

- ABC notation without metadata

- **Model**

- LSTM (long short-term memory)
- Working on the character level

*folk***RNN**
generate a folk tune with a recurrent neural network

PRESS TO GENERATE TUNE

Compose

MODEL
thesession.org (w/ :| |:)

TEMPERATURE SEED
1 62063

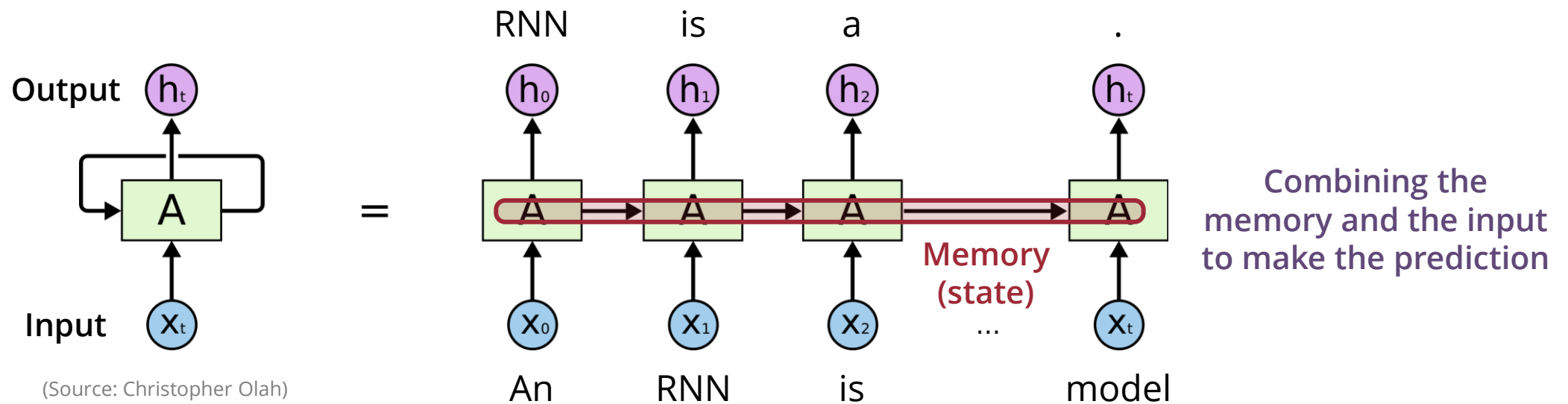
METER MODE
4/4 C Major

INITIAL ABC
Enter start of tune in ABC notation

folkrrnn.org

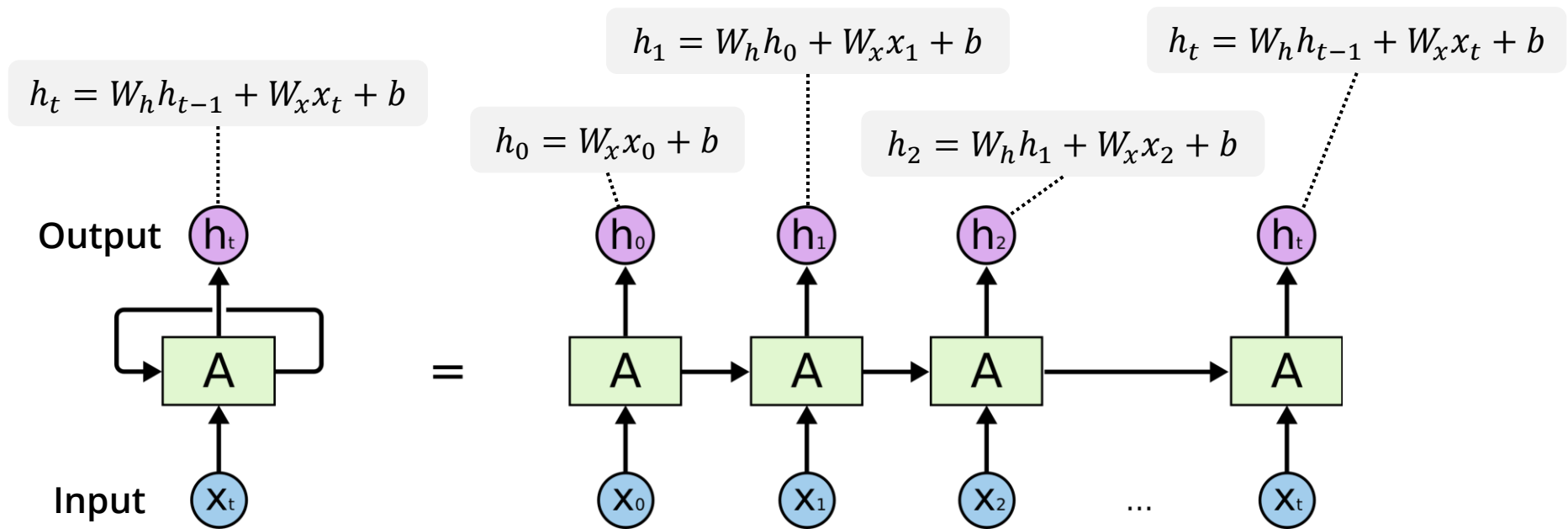
What is an RNN (Recurrent Neural Network)?

- A type of neural networks that have **loops**
- Widely used for **modeling sequences** (e.g., in natural language processing)



Vanilla RNNs

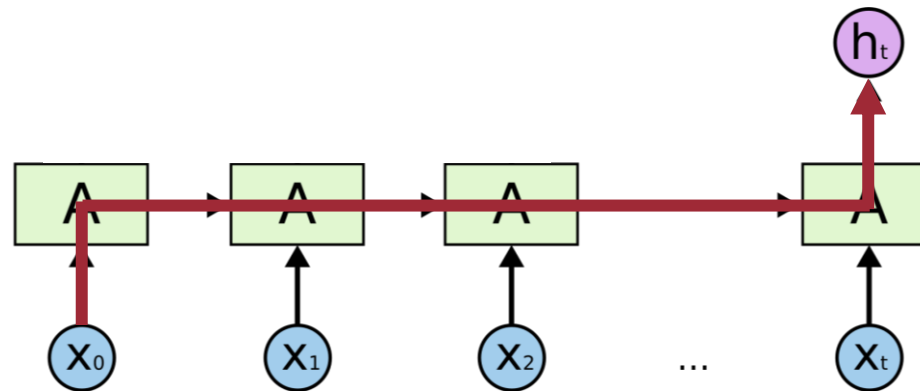
- The simplest form of RNNs
- LSTMs and GRUs are also RNNs



(Source: Christopher Olah)

Backpropagation Through Time

- An RNN is essentially a **very deep neural network**



$$h_t = W_h h_{t-1} + W_x x_t + b$$

$$h_t = W_h (W_h h_{t-2} + W_x x_{t-1} + b) + W_x x_t + b$$

⋮

$$h_t = W_h (W_x x_{t-1} + W_h (\dots W_h h_0 + W_x x_1 + b \dots) + b) + W_x x_t + b$$

Example: Folk RNN (Sturm et al., 2015)

- **Data**

- Collections of folk tunes

- **Representation**

- ABC notation without metadata

- **Model**

- **LSTM** (long short-term memory)
- Working on the **character** level

*folk***RNN**
generate a folk tune with a recurrent neural network

PRESS TO GENERATE TUNE

Compose

MODEL
thesession.org (w/ :| |:)

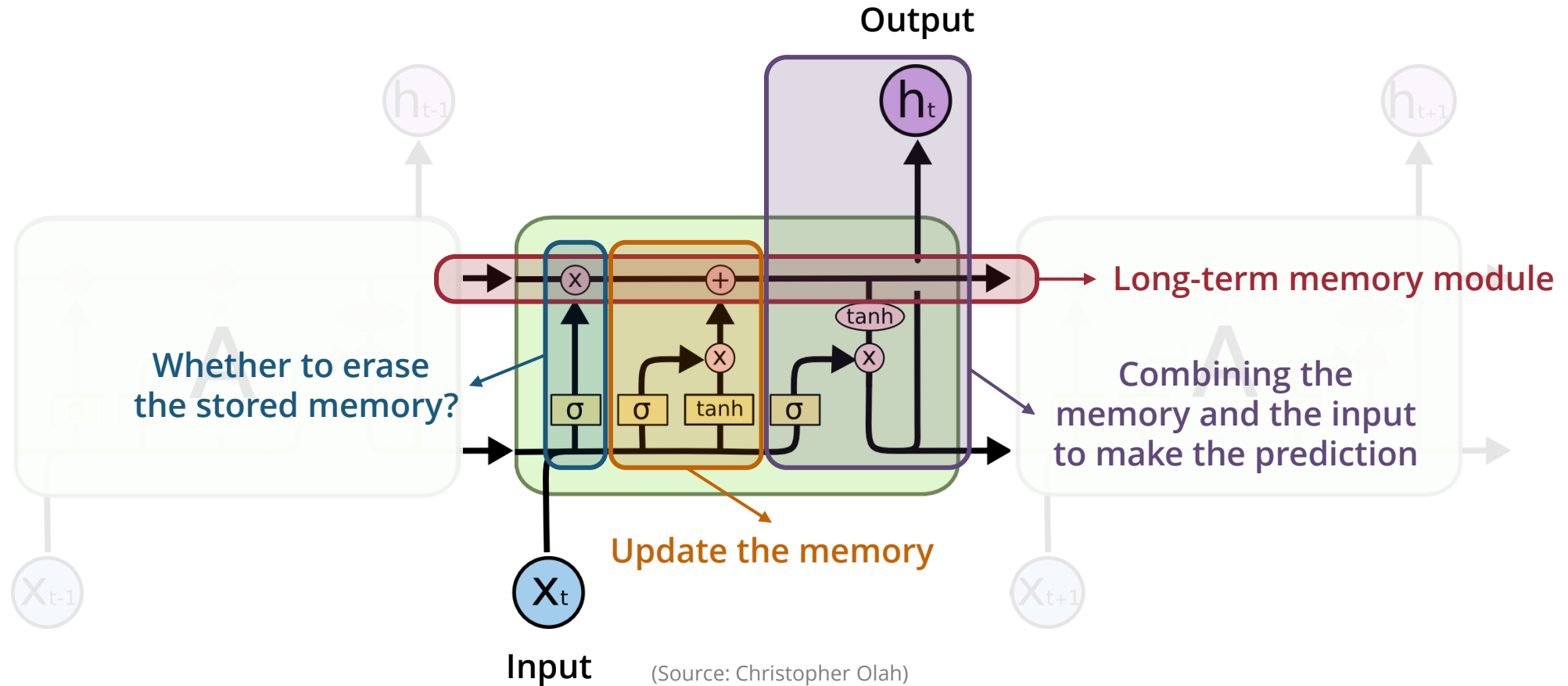
TEMPERATURE SEED
1 62063

METER MODE
4/4 C Major

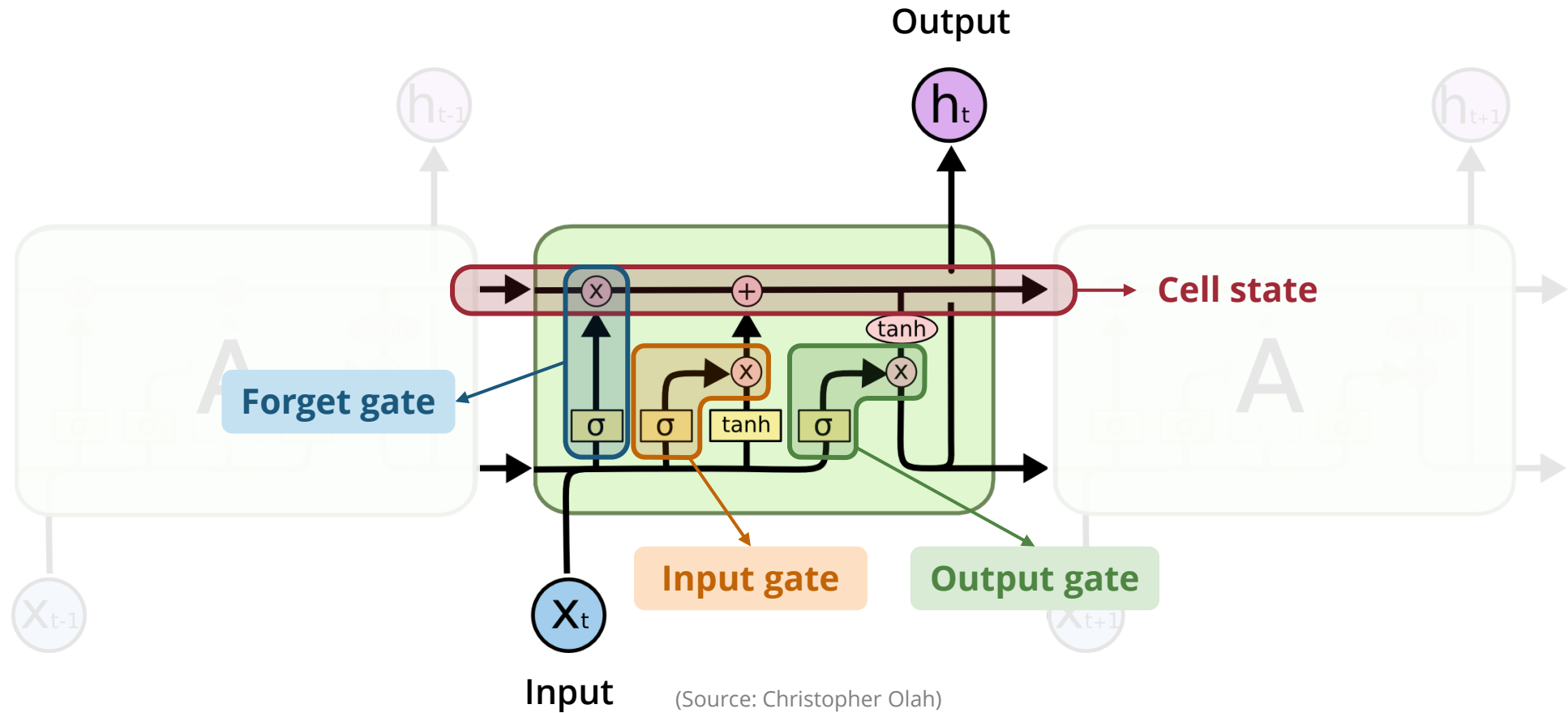
INITIAL ABC
Enter start of tune in ABC notation

folkrrnn.org

Demystifying LSTMs (Hochreiter & Schmidhuber, 1997)

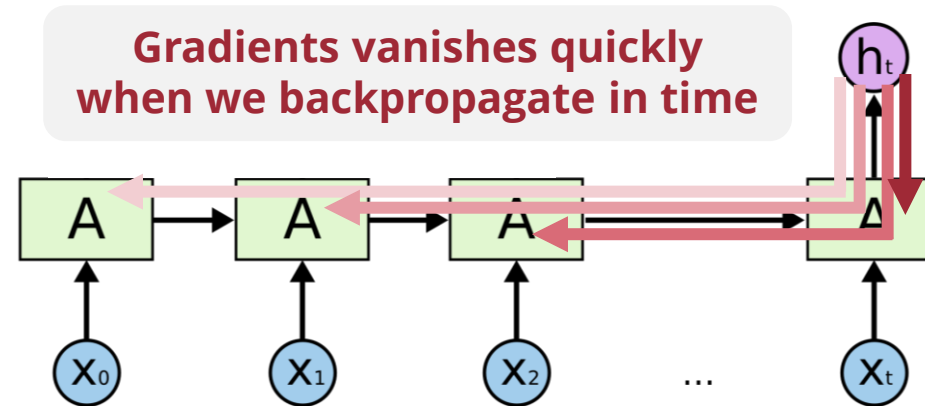


Demystifying LSTMs (Hochreiter & Schmidhuber, 1997)



Vanishing Gradients

- An RNN is essentially a **very deep neural network**



All the layers share the same weight matrix

Can still train the model without deeper gradients

Why bother?

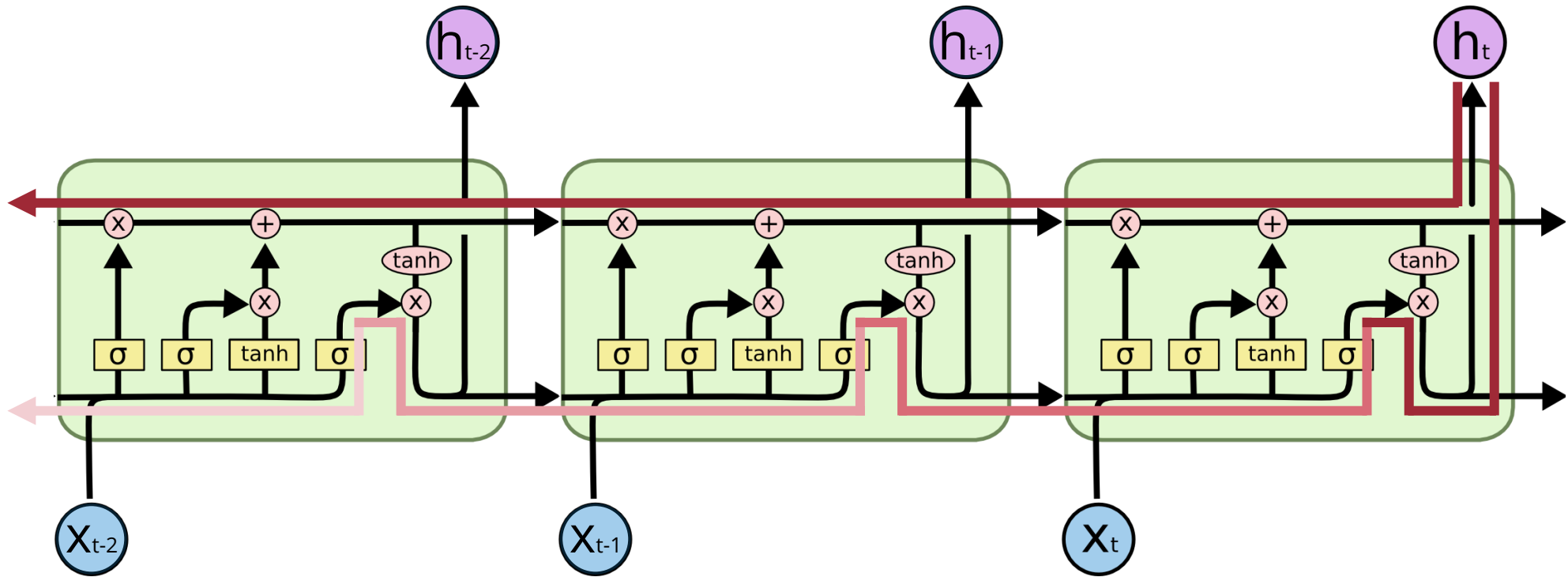
$$h_t = W_h h_{t-1} + W_x x_t + b$$

$$h_t = W_h (W_h h_{t-2} + W_x x_{t-1} + b) + W_x x_t + b$$

⋮

$$h_t = W_h (W_x x_{t-1} + W_h (\dots W_h h_0 + W_x x_1 + b \dots) + b) + W_x x_t + b$$

How can LSTMs Help Alleviate Vanishing Gradients?

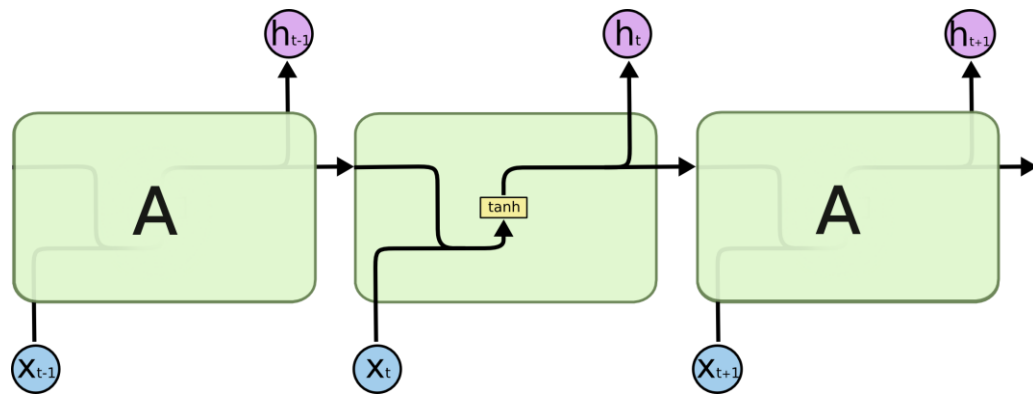


LSTMs does not completely solve vanishing gradients

Vanilla RNNs vs LSTMs

Vanilla RNN

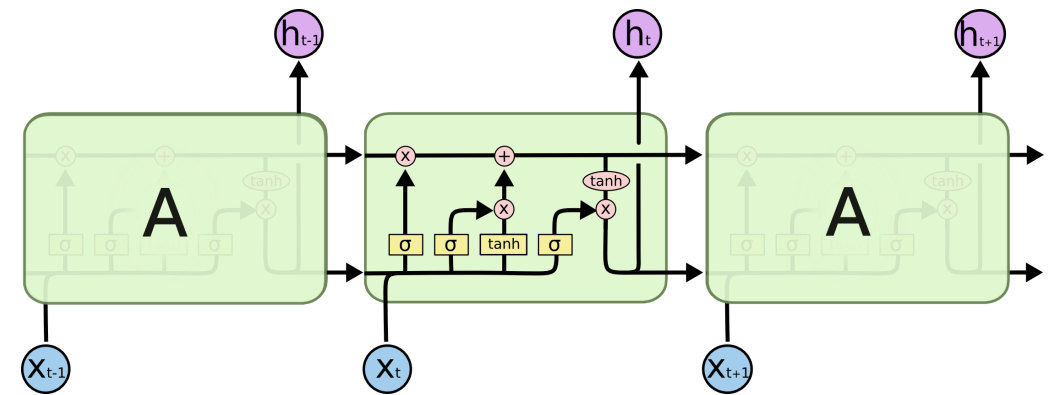
- Simplest form of RNNs
- Limited long-term memory
- Harder to train (due to gradient vanishing)



(Source: Christopher Olah)

LSTM

- Improved memory module
- Better long-term memory
- Easier to train



(Source: Christopher Olah)

Example: Folk RNN (Sturm et al., 2015)

- Data
 - Collections of folk tunes
- Representation
 - ABC notation without metadata
- Model
 - LSTM (long short-term memory)
 - Working on the **character level**

*folk***RNN**
generate a folk tune with a recurrent neural network

PRESS TO GENERATE TUNE

Compose

MODEL
thesession.org (w/ :| |:)

TEMPERATURE SEED
1 62063

METER MODE
4/4 C Major

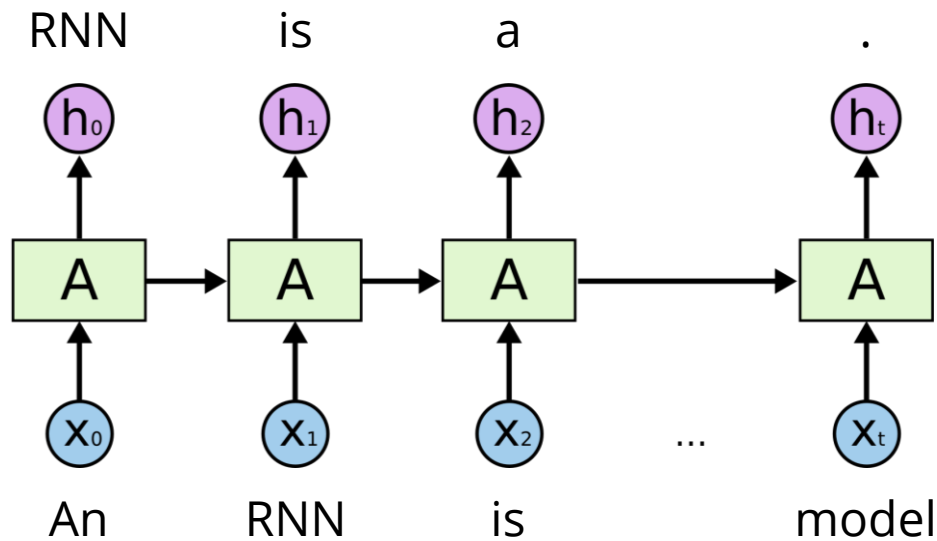
INITIAL ABC
Enter start of tune in ABC notation

folkrrnn.org

Word-level vs Character-level RNNs

Word-level RNNs

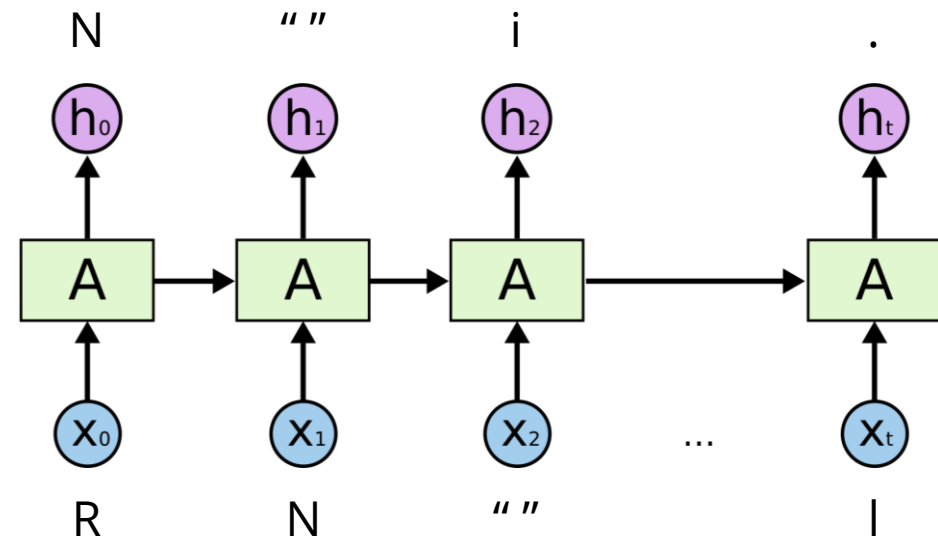
- Predicting word by word
- Most common



(Source: Christopher Olah)

Character-level RNNs

- Predicting character by character
- Useful when there is no natural "spaces"



(Source: Christopher Olah)

Limitations of ABC Notations

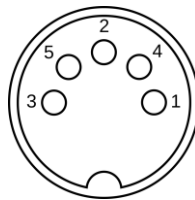
- Limited expressiveness
- Monophonic tunes only

MIDI-like Representation

MIDI (Musical Instrument Digital Interface)



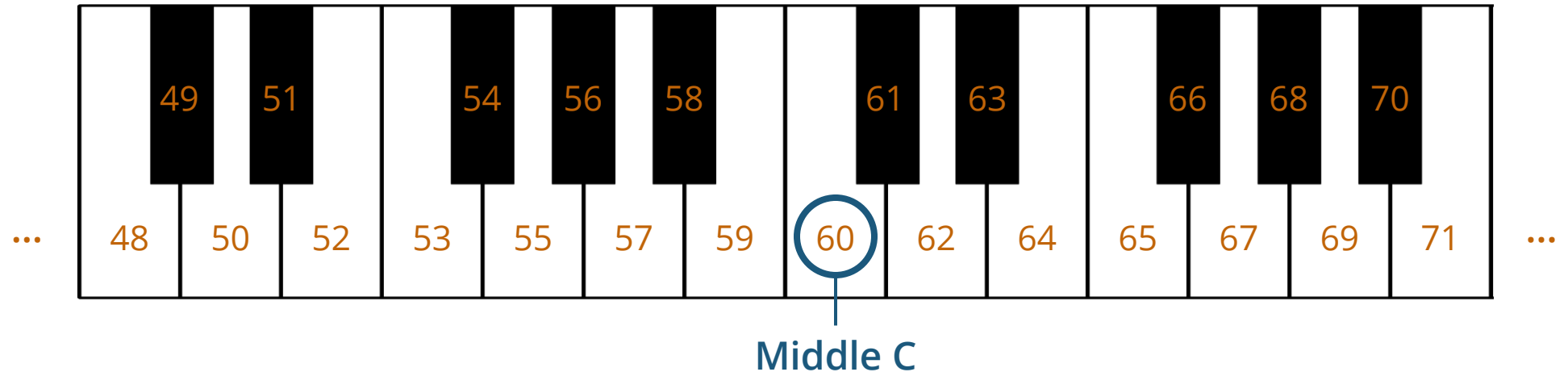
- A communication **protocol** between devices
- MIDI Messages
 - Note on
 - Note off
 - Delta time
 - Program change
 - Control change
 - Pitch bend change



MIDI I/O

MIDI Note Numbers


- Ranging from 0 to 127
 - Middle C is 60
 - Wider than standard piano's pitch range
- Widely used in various software, keyboards and algorithms



Representing Music using MIDI Messages

- Three main MIDI messages
 - Note on
 - Note off
 - Time Shift

Sunshine on the Meadow



The image shows two staves of musical notation. The top staff is in 4/4 time and contains a sequence of notes. The first note is circled in blue and has an orange 'X' over it. The second note is circled in green. A red arrow points from the first note to the second, indicating a time shift. The bottom staff continues the melody with more notes.

Note_on_67	Time_shift_quarter_note,	Note_off_67,
Note_on_67	Time_shift_quarter_note,	Note_off_67,
Note_on_64,	Time_shift_quarter_note,	Note_off_64,
Note_on_64,	Time_shift_quarter_note,	Note_off_64,
...		

Representing Polyphonic Music

- We can now handle music with multi-pitch at the same time
 - In the literature, “polyphonic” & “multi-pitch” are often used interchangeably

Clair de Lune
from “Suite Bergamasque” L. 75
3rd Movement
Claude Debussy
(1862–1918)

Andante très expressif

Piano

pp *con sordina*

Note_on_65, Note_on_68, Time_shift_eighth_note, Note_on_77, Note_on_80,
Time_shift_half_note, Note_off_77, Note_off_80, Note_on_73, Note_on_77,
Time_shift_dotted_quarter_note, Note_off_65, Note_off_68, ...

Example: Performance RNN (Oore et al., 2020)

- Data
 - Yamaha e-Piano Competition dataset (MAESTRO)
- Representation
 - 128 Note-On events
 - 128 Note-Off events
 - 125 Time-Shift events (8ms–1s)
 - 32 Set-Velocity events Handle dynamics
- Model
 - LSTM

Examples of generated music



Example: **A.I. Duet** (Mann et al., 2016)



youtu.be/OZE1bfPtvZo

experiments.withgoogle.com/ai/ai-duet/view/



Example: Music Transformer (Huang et al., 2019)

- **Data:** Yamaha e-Piano Competition dataset (MAESTRO)

- **Representation**

- 128 Note-On events
- 128 Note-Off events
- 100 Time-Shift events (10ms–1s)
- 32 Set-Velocity events

Almost the same representation as PerformanceRNN

Handle dynamics

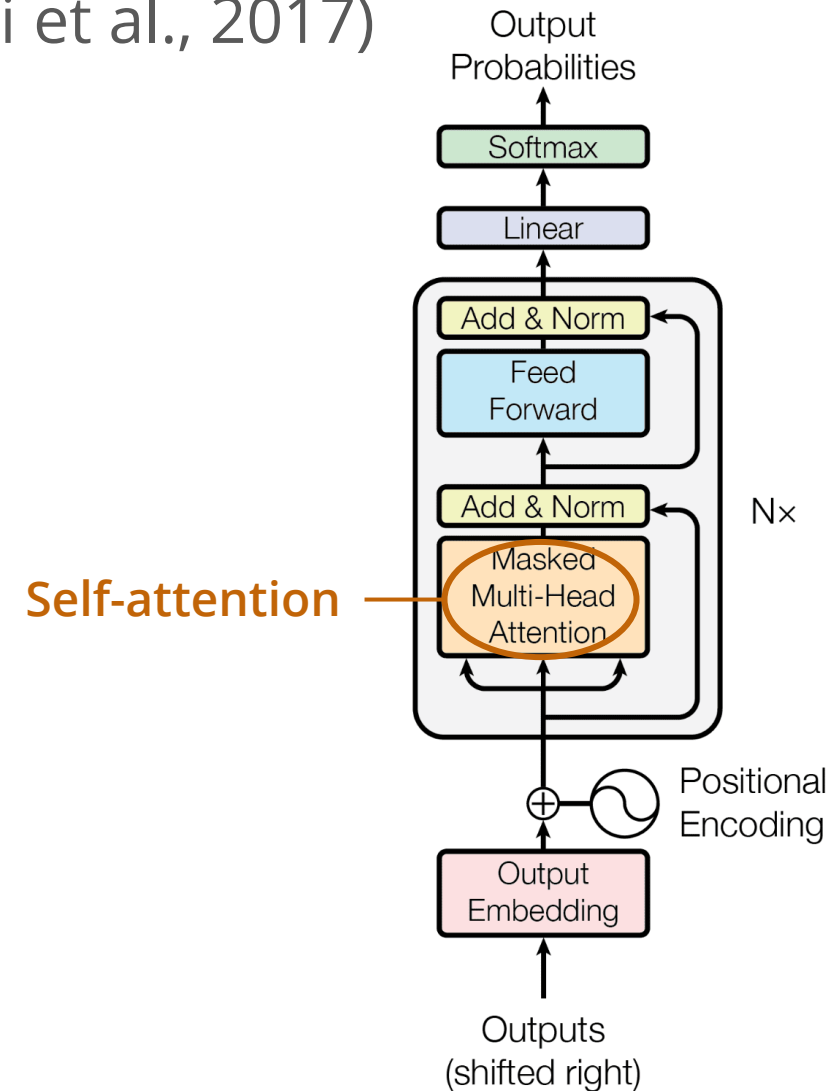
- **Model:** Transformer

Examples of generated music



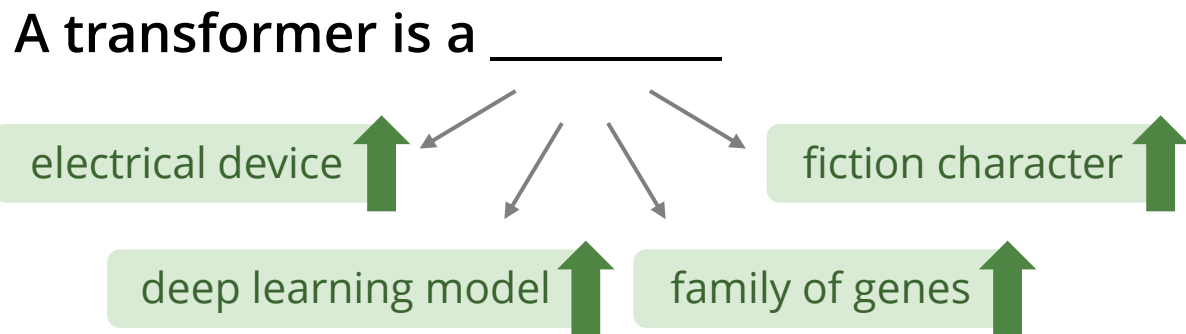
What is a Transformer? (Vaswani et al., 2017)

- A type of neural network that use the **self-attention mechanism**

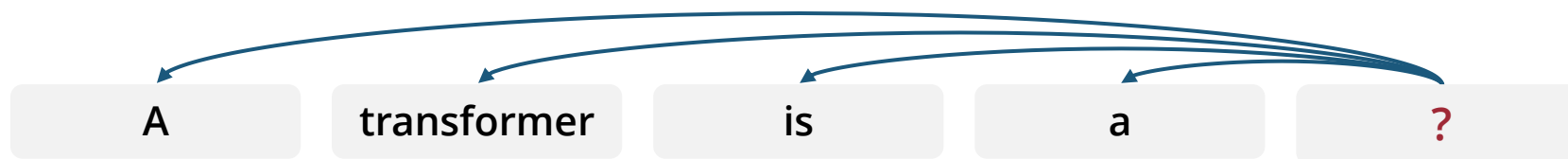


(Source: Vaswani et al., 2017; adapted)

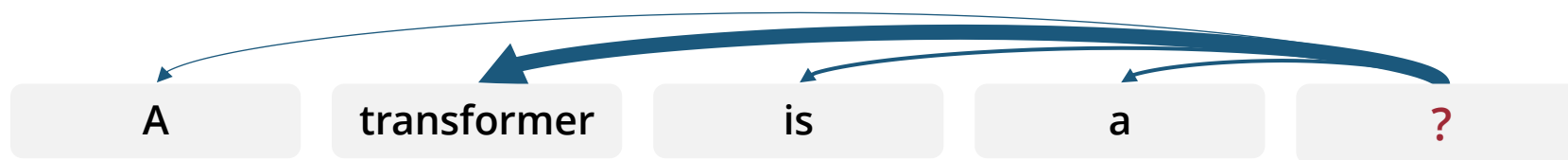
Self-attention Mechanism (Cheng et al., 2016)



Uniform attention

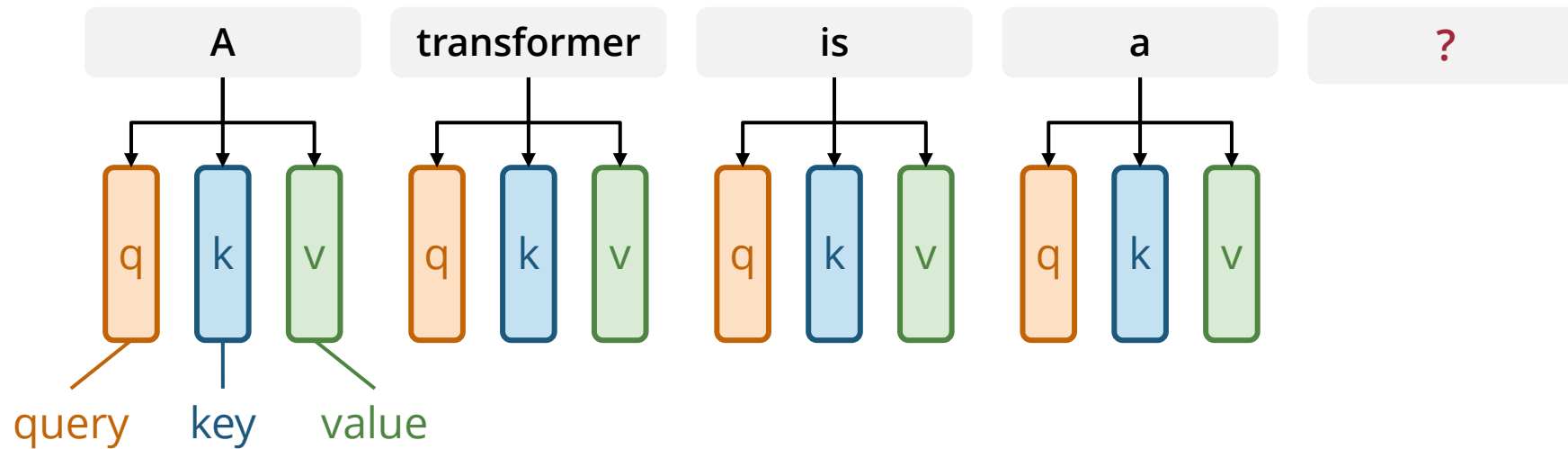


Variable attention

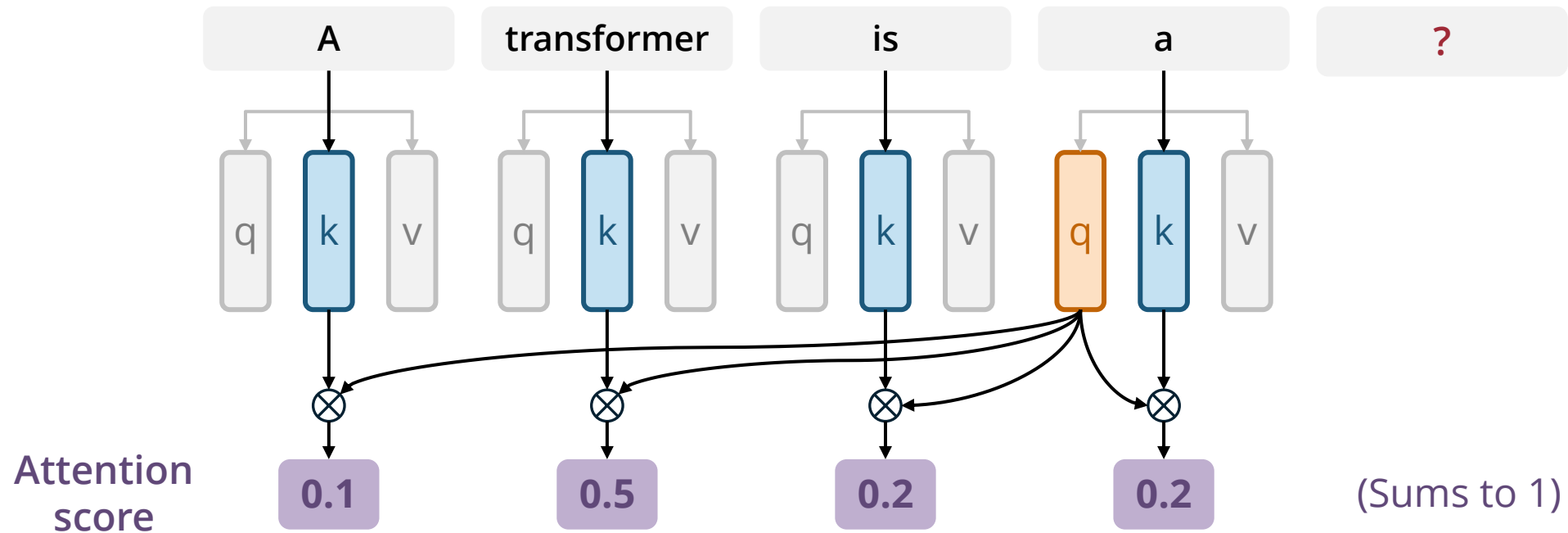


Transformers learn what to attend to from big data!

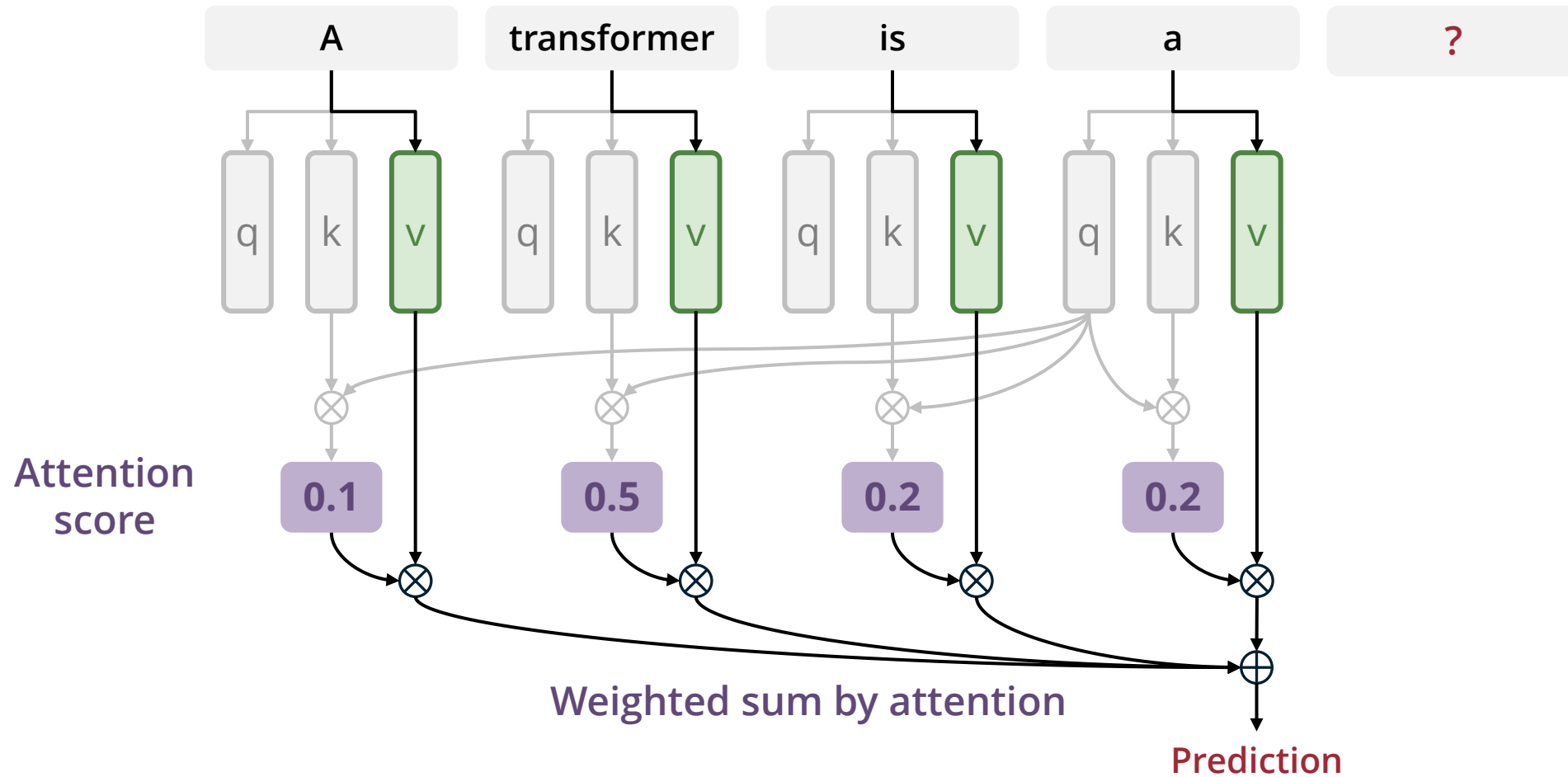
Demystifying Transformers (Vaswani et al., 2017)



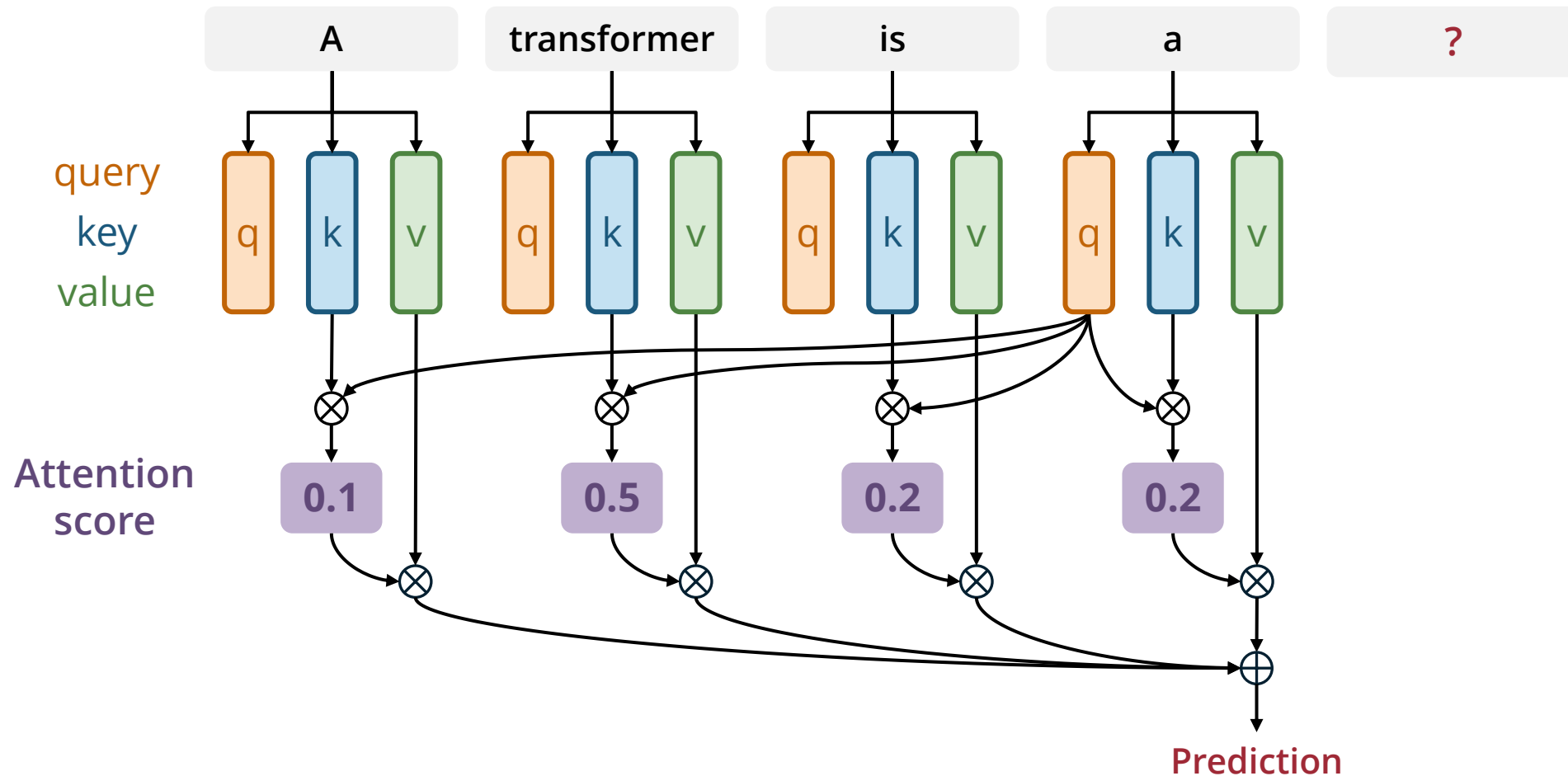
Demystifying Transformers (Vaswani et al., 2017)



Demystifying Transformers (Vaswani et al., 2017)



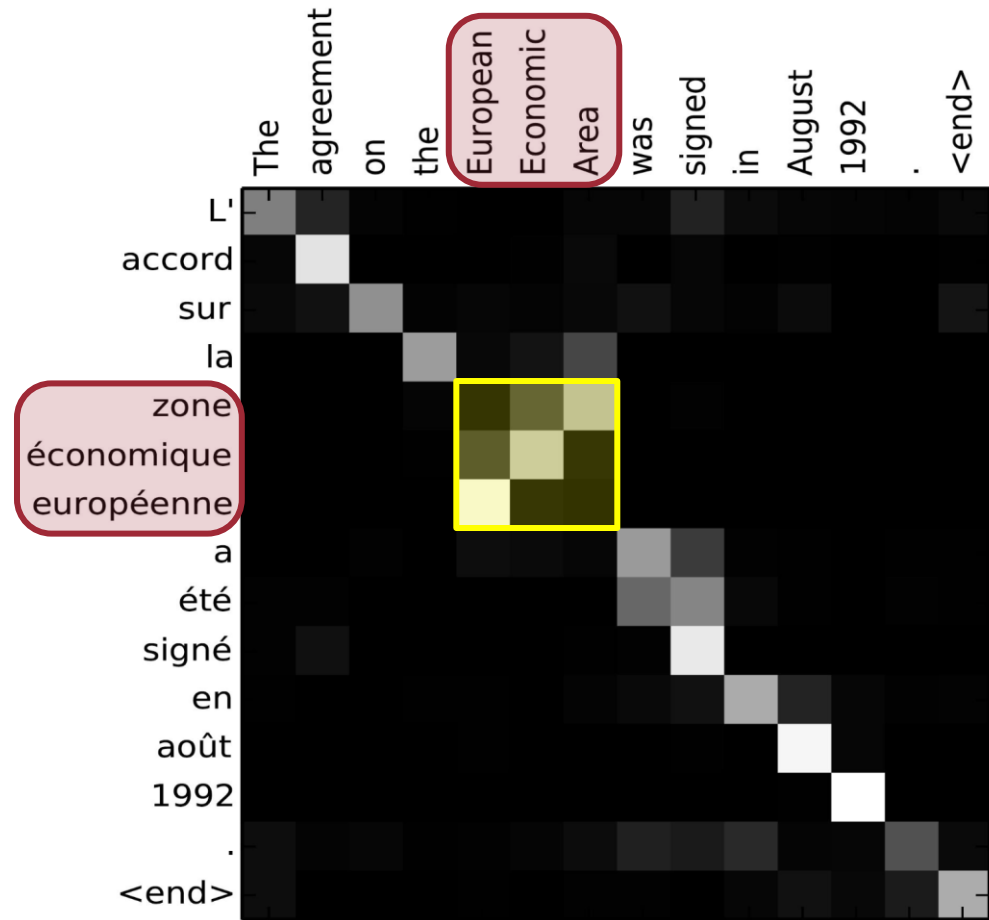
Demystifying Transformers (Vaswani et al., 2017)



Why Attention Mechanism?

The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .

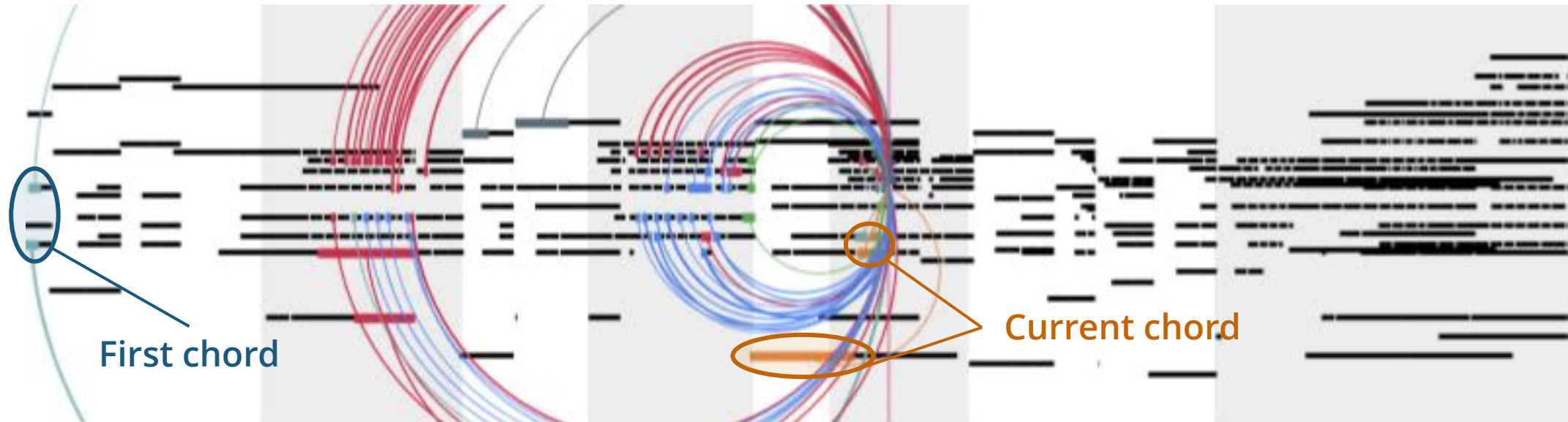
(Source: Cheng et al., 2016)



(Source: Bahdanau et al., 2015)

Visualizing Musical Self-attention

(Each color represents an attention head)



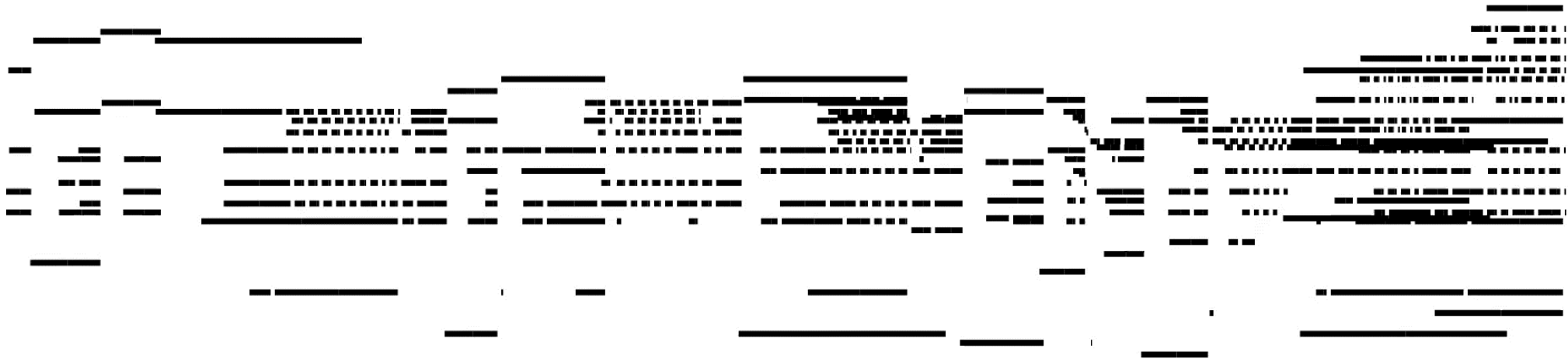
(Source: Huang et al., 2018)

Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Noam Shazeer, Ian Simon, Curtis Hawthorne, Andrew M. Dai, Matthew D. Hoffman, Monica Dinculescu, and Douglas Eck, "Music Transformer: Generating Music with Long-Term Structure," *ICLR*, 2019.

Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Noam Shazeer, Ian Simon, Curtis Hawthorne, Andrew M. Dai, Matthew D. Hoffman, Monica Dinculescu, and Douglas Eck, "Music Transformer: Generating Music with Long-Term Structure," *Magenta Blog*, December 13, 2018.

Visualizing Musical Self-attention

(Each color represents an attention head)



(Source: Huang et al., 2018)

Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Noam Shazeer, Ian Simon, Curtis Hawthorne, Andrew M. Dai, Matthew D. Hoffman, Monica Dinculescu, and Douglas Eck, "Music Transformer: Generating Music with Long-Term Structure," *ICLR*, 2019.

Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Noam Shazeer, Ian Simon, Curtis Hawthorne, Andrew M. Dai, Matthew D. Hoffman, Monica Dinculescu, and Douglas Eck, "Music Transformer: Generating Music with Long-Term Structure," *Magenta Blog*, December 13, 2018.

Beyond Solo Music

Representing Multiple Instruments

- Using **MIDI program change** messages

- Program numbers: 1–128 (or 0–127)
- 128 instruments in 16 families

Prog#	INSTRUMENT
1-8 PIANO	
1	Acoustic Grand
2	Bright Acoustic
3	Electric Grand
4	Honky-Tonk
5	Electric Piano 1
6	Electric Piano 2
7	Harpsichord
8	Clav

Prog#	INSTRUMENT	Prog#	INSTRUMENT	65-72 REED	73-80 PIPE		
1-8 PIANO		9-16 CHROMATIC PERCUSSION		65	Soprano Sax	73	Piccolo
1	Acoustic Grand	9	Celesta	66	Alto Sax	74	Flute
2	Bright Acoustic	10	Glockenspiel	67	Tenor Sax	75	Recorder
3	Electric Grand	11	Music Box	68	Baritone Sax	76	Pan Flute
4	Honky-Tonk	12	Vibraphone	69	Oboe	77	Blown Bottle
5	Electric Piano 1	13	Marimba	70	English Horn	78	Shakuhachi
6	Electric Piano 2	14	Xylophone	71	Bassoon	79	Whistle
7	Harpsichord	15	Tubular Bells	72	Clarinet	80	Ocarina
8	Clav	16	Dulcimer				
17-24 ORGAN		25-32 GUITAR		81-88 SYNTH LEAD		89-96 SYNTH PAD	
17	Drawbar Organ	25	Acoustic Guitar(nylon)	81	Lead 1 (square)	89	Pad 1 (new age)
18	Percussive Organ	26	Acoustic Guitar(steel)	82	Lead 2 (sawtooth)	90	Pad 2 (warm)
19	Rock Organ	27	Electric Guitar(jazz)	83	Lead 3 (calliope)	91	Pad 3 (polysynth)
20	Church Organ	28	Electric Guitar(clean)	84	Lead 4 (chiff)	92	Pad 4 (choir)
21	Reed Organ	29	Electric Guitar(muted)	85	Lead 5 (charang)	93	Pad 5 (bowed)
22	Accoridan	30	Overdriven Guitar	86	Lead 6 (voice)	94	Pad 6 (metallic)
23	Harmonica	31	Distortion Guitar	87	Lead 7 (fifths)	95	Pad 7 (halo)
24	Tango Accordion	32	Guitar Harmonics	88	Lead 8 (bass+lead)	96	Pad 8 (sweep)
33-40 BASS		41-48 STRINGS		97-104 SYNTH EFFECTS		105-112 ETHNIC	
33	Acoustic Bass	41	Violin	97	FX 1 (rain)	105	Sitar
34	Electric Bass(finger)	42	Viola	98	FX 2 (soundtrack)	106	Banjo
35	Electric Bass(pick)	43	Cello	99	FX 3 (crystal)	107	Shamisen
36	Fretless Bass	44	Cello	100	FX 4 (atmosphere)	108	Koto
37	Slap Bass 1	45	Contrabass	101	FX 5 (brightness)	109	Kalimba
38	Slap Bass 2	46	Tremolo Strings	102	FX 6 (goblins)	110	Bagpipe
39	Synth Bass 1	47	Pizzicato Strings	103	FX 7 (echoes)	111	Fiddle
40	Synth Bass 2	48	Orchestral strings	104	FX 8 (sci-fi)	112	Shanai
49-56 ENSEMBLE		57-64 BRASS		113-120 PERCUSSIVE		121-128 SOUND EFFECTS	
49	String Ensemble 1	57	Trumpet	113	Tinkle Bell	121	Guitar Fret Noise
50	String Ensemble 2	58	Trombone	114	Agogo	122	Breath Noise
51	SynthStrings 1	59	Tuba	115	Steel Drums	123	Seashore
52	SynthStrings 2	60	Muted Trumpet	116	Woodblock	124	Bird Tweet
53	Choir Aahs	61	French Horn	117	Telephone Ring	125	Telephone Ring
54	Voice Oohs	62	Brass Section	118	Melodic Tom	126	Helicopter
55	Synth Voice	63	SynthBrass 1	119	Synth Drum	127	Applause
56	Orchestra Hit	64	SynthBrass 2	120	Reverse Cymbal	128	Gunshot

(Source: Roger Dannenberg)

Example: MuseNet (Payne et al., 2019)

- **Data:** ClassicalArchives + BitMidi + MAESTRO
- **Representation:** “**instrument:velocity:pitch**”
 - Time shifts in real time (sec)
- **Model:** Transformer

```
bach piano_strings start tempo90
piano:v72:G1 piano:v72:G2 piano:v72:B4
piano:v72:D4 violin:v80:G4 piano:v72:G4
piano:v72:B5 piano:v72:D5 wait:12
piano:v0:B5 wait:5 piano:v72:D5 wait:12
...
```

Example of
generated music



Example: Multitrack Music Machine (Ens & Pasquier, 2020)

- **Data:** Lakh MIDI Dataset (LMD)
- **Representation:** as shown →
- **Model:** Transformer



LETS START WITH SOME U2

youtu.be/NdeMZ3y-84Q

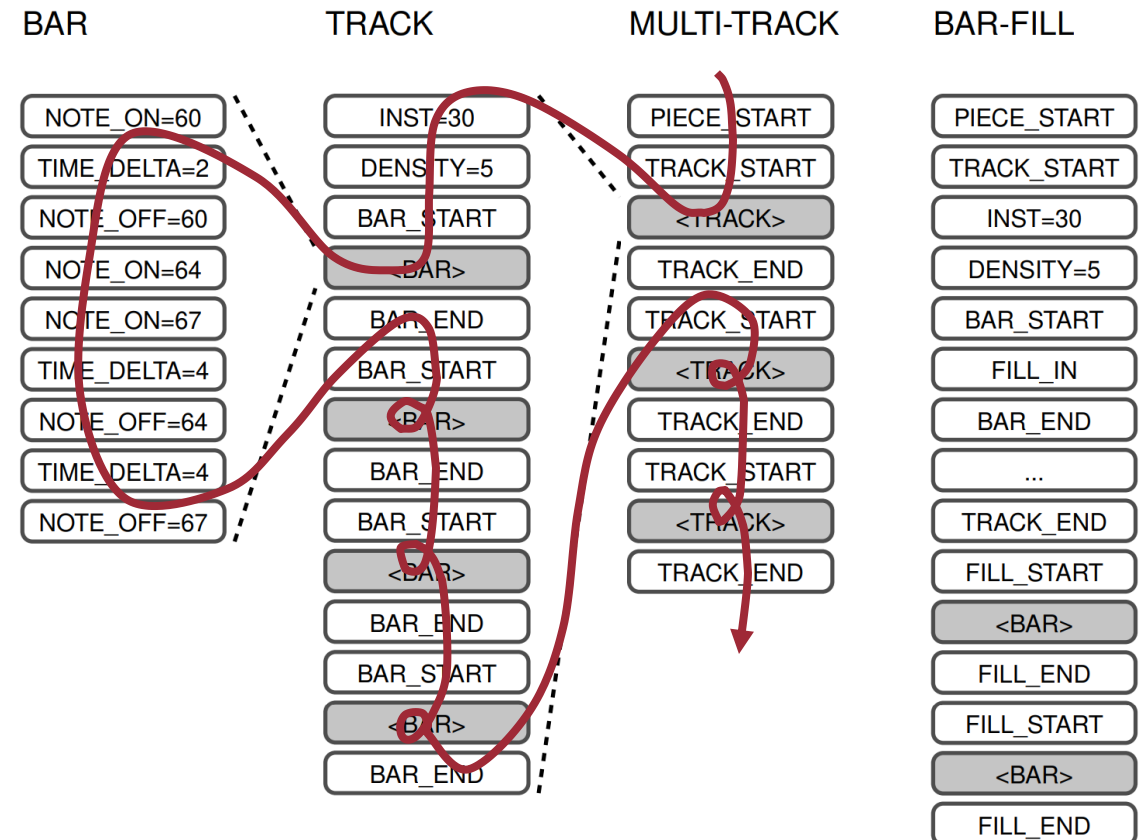


Fig. 1. The MultiTrack and BarFill representations are shown. The <bar> tokens correspond to complete bars, and the <track> tokens correspond to complete tracks.

(Ens & Pasquier, 2020)

Example: Multitrack Music Transformer (Dong et al., 2023)

- **Data:** Symbolic Orchestral Database (SOD)
- **Representation:** “(beat, position, pitch, duration, instrument)”
 - No time shift events **Why?**
- **Model:** Multi-dimensional Transformer

(0, 0, 0, 0, 0, 0)	Start of song
(1, 0, 0, 0, 0, 15)	Instrument: accordion
(1, 0, 0, 0, 0, 36)	Instrument: trombone
(1, 0, 0, 0, 0, 39)	Instrument: brasses
(2, 0, 0, 0, 0, 0)	Start of notes
(3, 1, 1, 41, 15, 36)	Note: beat=1, position=1, pitch=E2, duration=48, instrument=trombone
(3, 1, 1, 65, 4, 39)	Note: beat=1, position=1, pitch=E4, duration=12, instrument=brasses
(3, 1, 1, 65, 17, 15)	Note: beat=1, position=1, pitch=E4, duration=72, instrument=accordion
(3, 1, 1, 68, 4, 39)	Note: beat=1, position=1, pitch=G4, duration=12, instrument=brasses
(3, 1, 1, 68, 17, 15)	Note: beat=1, position=1, pitch=G4, duration=72, instrument=accordion
(3, 1, 1, 73, 17, 15)	Note: beat=1, position=1, pitch=C5, duration=72, instrument=accordion
(3, 1, 13, 68, 4, 39)	Note: beat=1, position=13, pitch=G4, duration=12, instrument=brasses
(3, 1, 13, 73, 4, 39)	Note: beat=1, position=13, pitch=C5, duration=12, instrument=brasses
(3, 2, 1, 73, 12, 39)	Note: beat=2, position=1, pitch=C5, duration=36, instrument=brasses
(3, 2, 1, 77, 12, 39)	Note: beat=2, position=1, pitch=E5, duration=36, instrument=brasses
...	...
(4, 0, 0, 0, 0, 0)	End of song

(Source: Dong et al., 2023)

Example of
generated music



Drums in MIDI

- **Channel 10** is reserved for drums
- Encoded by MIDI pitches 35–81
- Models that support drums
 - **MuseNet** (Payne et al., 2019)
 - **Song from PI** (Chu et al., 2017)
 - **MMM** (Ens and Pasquier, 2019)
 - *and many more...*

	A2		
Acoustic Bass Drum (35)	B2	A#2	
Bass Drum 1 (36)	C3		
Acoustic Snare (38)	D3	C#3	(37) Side Stick
Electric Snare (40)	E3	D#3	(39) Hand Clap
Low Floor Tom (41)	F3		
High Floor Tom (43)	G3	F#3	(42) Closed Hi-Hat
Low Tom (45)	A3	G#3	(44) Pedal Hi-Hat
Low-Mid Tom (47)	B3	A#3	(46) Open Hi-Hat
Hi-Mid Tom (48)	C4		
High Tom (50)	D4	C#4	(49) Crash Cymbal 1
Chinese Cymbal (52)	E4	D#4	(51) Ride Cymbal 1
Ride Bell (53)	F4		
Splash Cymbal (55)	G4	F#4	(54) Tambourine
Crash Cymbal 2 (57)	A4	G#4	(56) Cowbell
Ride Cymbal 2 (59)	B4	A#4	(58) Vibraslap
Hi Bongo (60)	C5		
Mute Hi Conga (62)	D5	C#5	(61) Low Bongo
Low Conga (64)	E5	D#5	(63) Open Hi Conga
High Timbale (65)	F5		
High Agogo (67)	G5	F#5	(66) Low Timbale
Cabasa (69)	A5	G#5	(68) Low Agogo
Short Whistle (71)	B5	A#5	(70) Maracas
Long Whistle (72)	C6		
Long Guiro (74)	D6	C#6	(73) Short Guiro
Hi Wood Block (76)	E6	D#6	(75) Claves
Low Wood Block (77)	F6		
Open Cuica (79)	G6	F#6	(78) Mute Cuica
Open Triangle (81)	A6	G#6	(80) Mute Triangle

(Source: Wikipedia)

en.wikipedia.org/wiki/General_MIDI

Christine Payne, "MuseNet," *OpenAI*, 2019.

Hang Chu, Raquel Urtasun, and Sanja Fidler, "Song From PI: A Musically Plausible Network for Pop Music Generation," *ICLR Workshop*, 2017.

Jeff Ens and Philippe Pasquier, "MMM : Exploring Conditional Multi-Track Music Generation with the Transformer," *arXiv preprint arXiv:2008.06048*, 2020.

The Many Representations for Music Generation

- **PerformanceRNN** (Oore et al., 2020)
- **REMI** (Huang et al., 2020)
- **MuMIDI** (Ren et al., 2020)
- **Compound Word** (Hsiao et al., 2021)
- **REMI+** (von Rütte et al., 2023)
- **TSD** (Fradet et al., 2023)
- *and so on...*

MIDITok

github.com/Natooz/MidiTok



Sageev Oore, Ian Simon, Sander Dieleman, Douglas Eck, and Karen Simonyan, "This Time with Feeling: Learning Expressive Musical Performance", *Neural Computing and Applications*, 32, 2020.

Yu-Siang Huang and Yi-Hsuan Yang, "Pop Music Transformer: Beat-based Modeling and Generation of Expressive Pop Piano Compositions," *MM*, 2020.

Yi Ren, Jinzheng He, Xu Tan, Tao Qin, Zhou Zhao, and Tie-Yan Liu, "PopMAG: Pop Music Accompaniment Generation," *MM*, 2020.

Wen-Yi Hsiao, Jen-Yu Liu, Yin-Cheng Yeh, and Yi-Hsuan Yang, "Compound Word Transformer: Learning to Compose Full-Song Music over Dynamic Directed Hypergraphs," *AAAI*, 2021.

Dimitri von Rütte, Luca Biggio, Yannic Kilcher, and Thomas Hofmann, "FIGARO: Generating Symbolic Music with Fine-Grained Artistic Control," *ICLR*, 2023.

Nathan Fradet, Nicolas Gutowski, Fabien Chhel, and Jean-Pierre Briot, "Byte Pair Encoding for Symbolic Music," *EMNLP*, 2023.

Symbolic Music Datasets

- [JSBach Chorale](#)
- [MusicNet](#)
- [Essen Folk Song Dataset](#)
- [Wikifonia](#)
- [Lakh MIDI Dataset](#)
- [MetaMIDI](#)
- Expressive MIDI: [MAESTRO](#)

Symbolic Music Datasets

Dataset	Format	Hours	Songs	Genre
Lakh MIDI Dataset	MIDI	>5000	174,533	misc
MAESTRO Dataset	MIDI	201.21	1,282	classical
Wikifonia Lead Sheet Dataset	MusicXML	198.40	6,405	misc
Essen Folk Song Dataset	ABC	56.62	9,034	folk
NES Music Database	MIDI	46.11	5,278	game
MusicNet Dataset	MIDI	30.36	323	classical
Hymnal Tune Dataset	MIDI	18.74	1,756	hymn
Hymnal Dataset	MIDI	17.50	1,723	hymn
music21's Corpus	misc	16.86	613	misc
EMOPIA Dataset	MIDI	10.98	387	pop
Nottingham Database	ABC	10.54	1,036	folk
music21's JSBach Corpus	MusicXML	3.46	410	classical
JSBach Chorale Dataset	MIDI	3.21	382	classical
Haydn Op.20 Dataset	Humdrum	1.26	24	classical

(Source: MusPy Documentation)