

2018 Fall  
**CTP431: Music and Audio Computing**

# Automatic Music Generation

Graduate School of Culture Technology, KAIST  
Juhan Nam



# Outlines

- Early Approaches
  - Markov Models
  - Recombinant Models
  - Cellular Automata
  - Genetic Algorithm
- Recent Advances
  - Neural Networks
- Interactive music generation



# Symbolic Music

- Symbolic music is represented as a sequence of notes

Sonate No. 8, “Pathétique”

3rd Movement  
Opus 13

Ludwig van Beethoven  
(1770 - 1827)

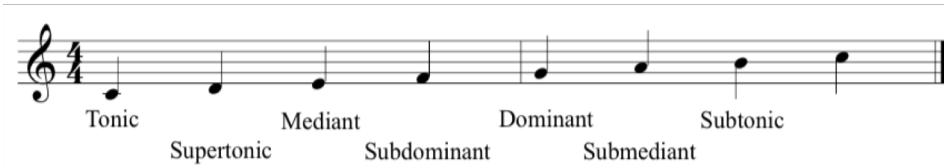
Piano

Rondo Allegro

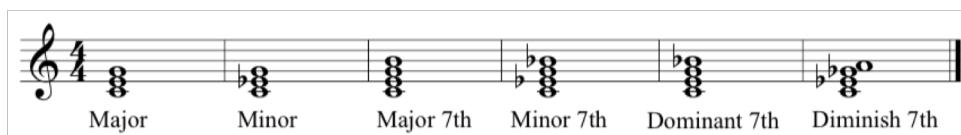
The musical score consists of two staves of piano music. The top staff begins with a dynamic marking 'p' (pianissimo). The bottom staff starts with a measure of eighth-note pairs. The score is framed by a blue border.

# Symbolic Music

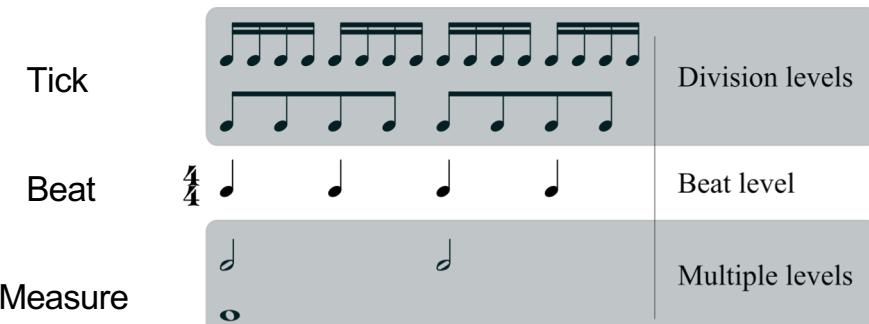
- Music is structured sequential data



Scale



Harmony



Rhythm



Form

# Symbolic Music

- Musical notes are temporally dependent
  - Note-level
  - Beat-level
  - Measure-level



# Markov Model

- A random variable  $q$  has  $N$  states ( $S_1, S_2, \dots, S_N$ ) and, at each time step, one of the states are randomly chosen:  $q_t \in \{S_1, S_2, \dots, S_N\}$
- The probability distribution for the current state is determined by the previous state(s)
  - The first-order Markov model:  $P(q_t|q_1, q_2, \dots, q_{t-1}) = P(q_t|q_{t-1})$
  - The second-order Markov model:  $P(q_t|q_1, q_2, \dots, q_{t-1}) = P(q_t|q_{t-1}, q_{t-2})$

# Markov Model

- Example: simple melody generation

- $q_t \in \{C, D, E\}$
- The transition probability matrix 3 by 3

$$P(q_t = C | q_{t-1} = C) = 0.7$$

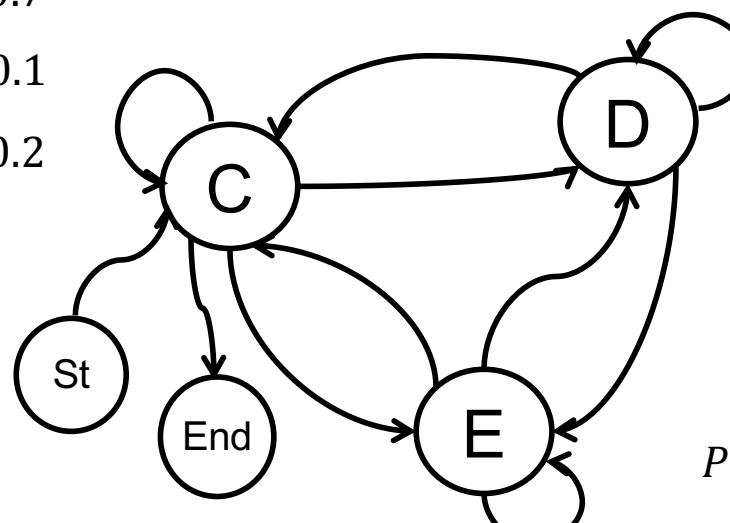
$$P(q_t = D | q_{t-1} = C) = 0.1$$

$$P(q_t = E | q_{t-1} = C) = 0.2$$

$$P(q_t = C | q_{t-1} = D) = 0.2$$

$$P(q_t = D | q_{t-1} = D) = 0.6$$

$$P(q_t = E | q_{t-1} = D) = 0.2$$



$$P(q_t = C | q_{t-1} = E) = 0.3$$

$$P(q_t = D | q_{t-1} = E) = 0.1$$

$$P(q_t = E | q_{t-1} = E) = 0.6$$

# Markov Model

- The transition matrix can be learned from data
  - Dancing Markov Gymnopédies: <https://codepen.io/teropa/pen/bRqYVj/>
- Generated music
  - Learned with Satie's "Gymnopédies" and "Trois Gnossiennes"
    - <https://www.youtube.com/watch?v=H3xgdDTvvlc>
  - Learned with Bach's "Toccata and Fugue in D minor" (BWV 565)
    - <https://www.youtube.com/watch?v=lOIaK0x4vA>

# Example: Illiac Suite

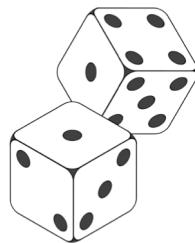
- The first computer-generated composition (1956)
  - Lejaren Hiller and Leonard Issacson
  - They used Markov models of variable order to select notes with different lengths
- Music
  - <https://www.youtube.com/watch?v=n0njBFLQSk8&list=PLIVblwUBdcStsNpl0v4OCbC5k-mIDcyaR>



# Recombinant Music

- Musical Dice Game

- Generate from pre-composed small pieces by random draws
- The table of me preserves musical “style”



|                  |  | A  | B   | C   | D   | E   | F   | G   | H   |     |
|------------------|--|----|-----|-----|-----|-----|-----|-----|-----|-----|
| Erster Theil.    |  | 12 | 96  | 92  | 141 | +1  | 105 | 122 | 11  | 30  |
| Premiere Partie. |  | 3  | 134 | 6'  | 128 | 63  | 140 | 46  | 134 | 81  |
| Zweiter Theil.   |  | 4  | 69  | 93  | 138 | 13  | 133 | 35  | 110 | 24  |
| Seconde Partie.  |  | 5  | 40  | 17  | 117 | 85  | 161 | 2   | 159 | 100 |
|                  |  | 6  | 148 | 74  | 163 | 43  | 80  | 37  | 36  | 107 |
|                  |  | 7  | 104 | 137 | 27  | 167 | 15+ | 6%  | 118 | 91  |
|                  |  | 8  | 169 | 60  | 171 | 33  | 99  | 133 | 21  | 127 |
|                  |  | 9  | 119 | 5+  | 114 | 30  | 140 | 86  | 169 | 94  |
|                  |  | 10 | 88  | 142 | 42  | 116 | 75  | 129 | 69  | 123 |
|                  |  | 11 | 3   | 87  | 165 | 61  | 133 | 47  | 147 | 93  |
|                  |  | 12 | 4+  | 130 | 10  | 103 | 28  | 37  | 106 | 4   |

A musical score titled "TABLE de MUSIQUE." featuring 32 numbered measures of music. The score is written in two systems of four staves each, with measures numbered 1 through 32 below each staff. The music consists of eighth and sixteenth note patterns with various dynamics and rests.

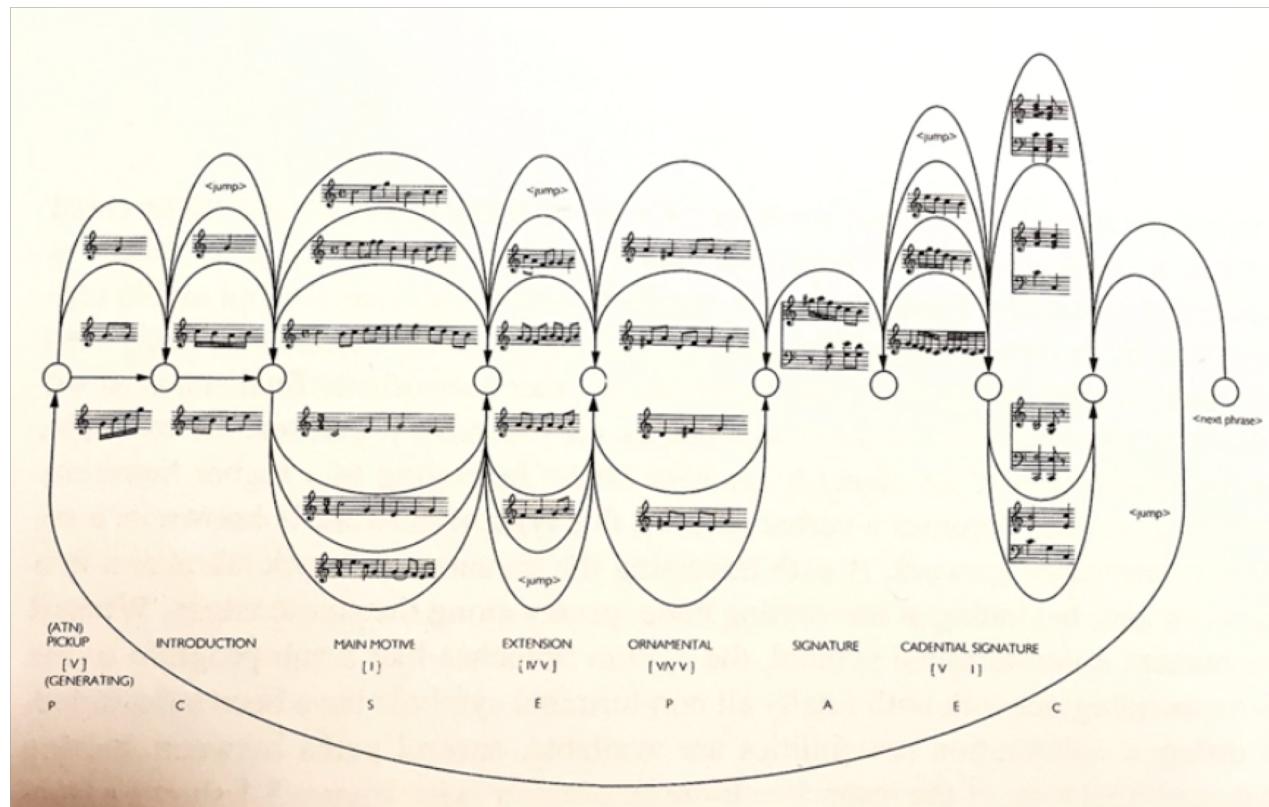
$$11^{16} = 45,949,729,863,572,161 \text{ variations}$$

[https://imslp.org/wiki/Musikalisches\\_W%BCrfelspiel,\\_K.516f\\_\(Mozart,\\_Wolfgang\\_Amadeus\)](https://imslp.org/wiki/Musikalisches_W%BCrfelspiel,_K.516f_(Mozart,_Wolfgang_Amadeus))

Mozart K. 516F

# Recombinant Music

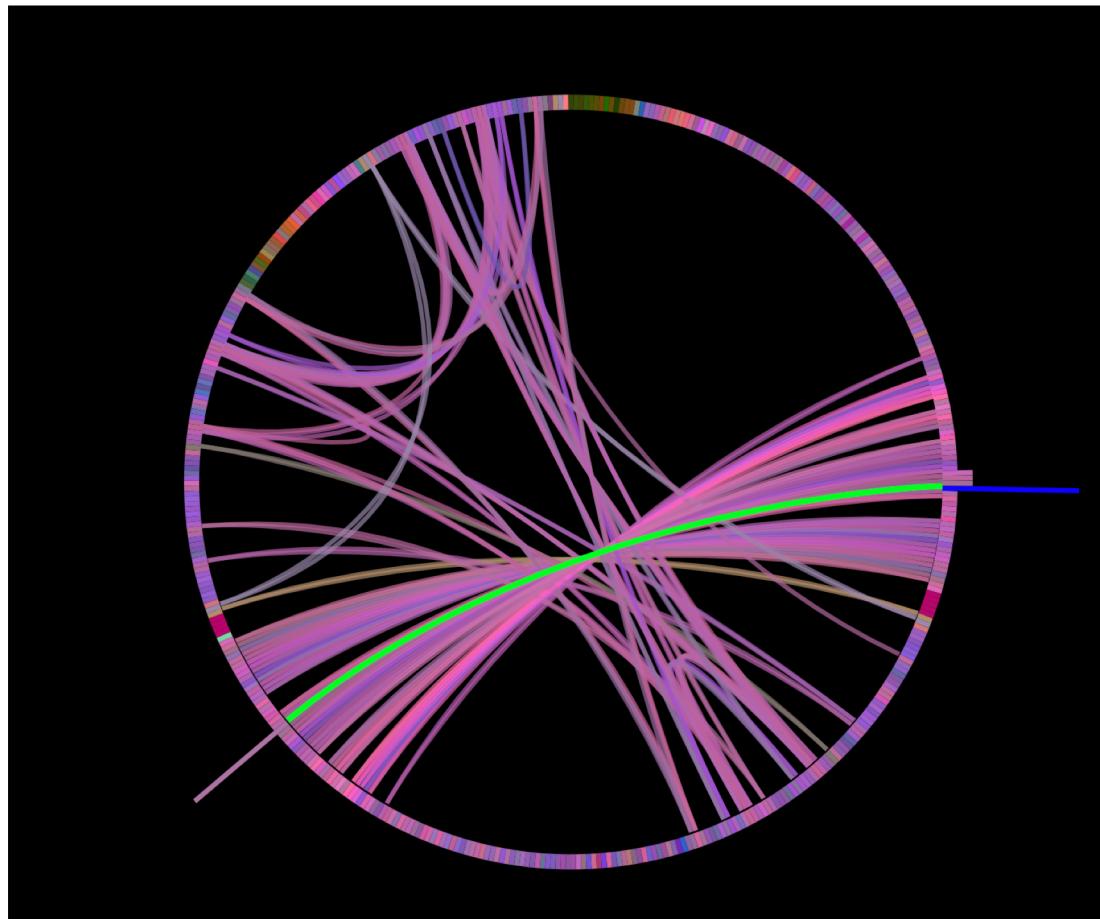
- David Cope's Experiments in Musical Intelligence (EMI)
  - Segment and reassemble existing pieces of music by pattern matching
  - Create a new piece of music that preserves the style of the original



Augmented Transition Networks (David Cope)

# Infinite Jukebox

- Music mash-up using beat-level self-similarity within a song



<http://infinitejukebox.playlistmachinery.com/>

# “In C”

- Ted Riley’s ensemble music
  - Also called “Minimal music”

“In C”  
by Terry Riley

## Instruction for beginners

1 Any number of people can play this piece on any instrument or instruments (including voice).

2 The piece consists of 53 melodic patterns to be repeated any amount of times. You can choose to start a new pattern at any point. The choice is up to the individual performer! We suggest beginners are very familiar with patterns 1-12.

3 Performers move through the melodic patterns in order and cannot go back to an earlier pattern. Players should try to stay within 2-3 patterns of each other.

4 If any pattern is too technically difficult, feel free to move to the next one.

5 The eighth note pulse is constant. Always listen for this pulse. The pulse for our experience will be piano and Orff instruments being played on the stage.

6 The piece works best when all the players are listening very carefully. Sometimes it is better to just listen and not play. It is important to fit into the group sound and understand how what you decide to play affects everybody around you. If you play softly, other players might follow you and play soft. If you play loud, you might influence other players to play loud.

7 The piece ends when the group decides it ends. When you reach the final pattern, repeat it until the entire group arrives on this figure. Once everyone has arrived, let the music slowly die away.

The score consists of 53 numbered melodic patterns, each on a single-line staff. Patterns 1-12 are relatively simple, featuring mostly quarter notes and eighth notes. Patterns 13-20 introduce more complex rhythms like sixteenth-note figures. Patterns 21-30 continue this trend with increasing complexity. Patterns 31-40 feature sustained notes and grace notes. Patterns 41-53 show a variety of rhythmic patterns, including eighth-note chords and sixteenth-note runs. The patterns are separated by vertical bar lines and some horizontal measures.

Figure 1.1. Score of *In C* (copyright Terry Riley, 1964).

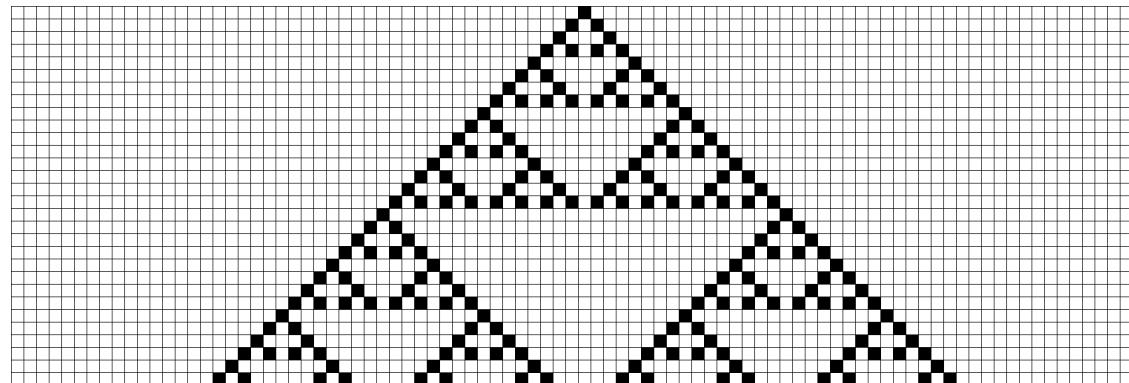
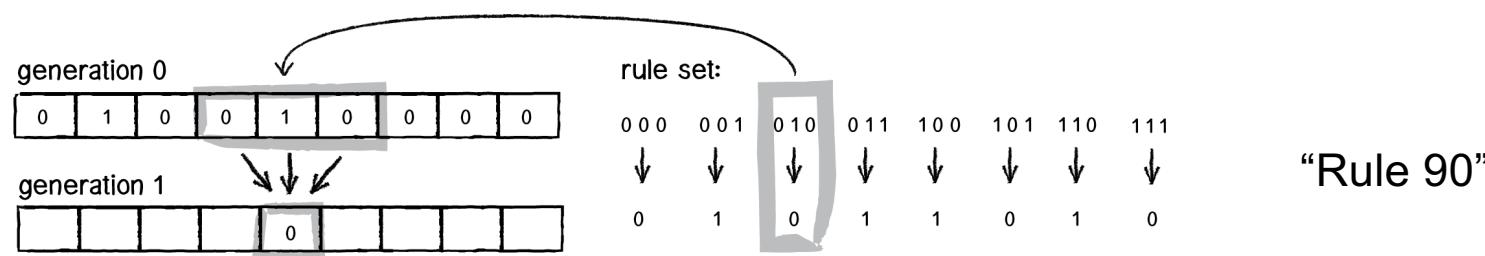
Source: <https://www.musicinst.org/sites/default/files/attachments/ln%20C%20Instructions%20for%20Beginners.pdf>

Source: <https://nmbx.newmusicusa.org/terry-rileys-in-c/>

[https://www.youtube.com/results?search\\_query=Terry+Riley+In+C](https://www.youtube.com/results?search_query=Terry+Riley+In+C)

# Cellular Automata

- A cell-based state evolution model
  - Determines the **state** of each **cell** using **neighbors** and **a rule set**
  - A Wolfram model example:



- Related to self-replicating patterns in biology

Source: <https://natureofcode.com/book/chapter-7-cellular-automata/>

# Conway's Game of Life

- 2D cellular automata

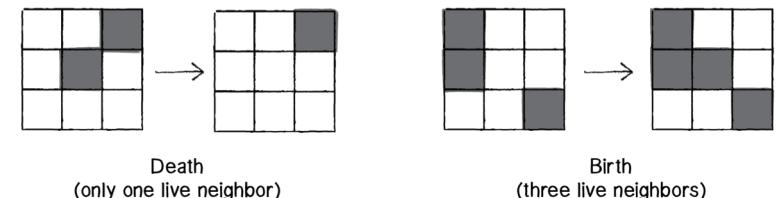
- Rules of life

- Death ( $1 \rightarrow 0$ ) : overpopulation ( $\geq 4$ ) or loneliness ( $\leq 1$ )
  - Birth ( $0 \rightarrow 1$ ) : 3 neighbors are alive
  - Otherwise, stay in the same state

Two-dimensional cellular automata

|   |   |   |   |   |   |
|---|---|---|---|---|---|
| 1 | 0 | 1 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 1 |
| 1 | 1 | 1 | 0 | 1 | 1 |
| 1 | 0 | 1 | 0 | 1 | 0 |
| 0 | 0 | 0 | 1 | 1 | 0 |
| 1 | 1 | 0 | 0 | 1 | 0 |
| 1 | 1 | 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | 1 | 1 | 1 |

a neighborhood of 9 cells



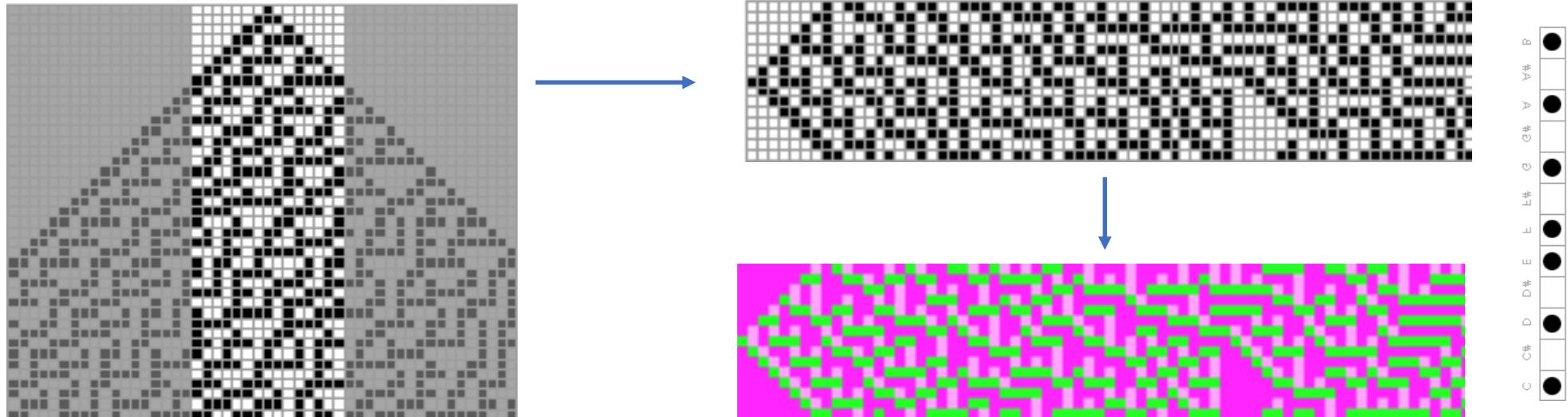
Source: <https://natureofcode.com/book/chapter-7-cellular-automata/>

- Demos:

- <http://www.cappel-nord.de/webaudio/conways-melodies/>
  - <http://nexusosc.com/gameofreich/>
  - <http://blipsoflife.herokuapp.com/>

# WolframTones

- Automatic music generation system based on cellular automata



Mapping to musical notes by rules

- Demo: <http://tones.wolfram.com/generate>

# Statistical Models

- As aforementioned, music is highly structured sequence data. Thus, we can model the sequence using **an auto-regressive model.**



$$p(q_t | q_1, \dots, q_{t-1})$$

$q_t$ : note features

- In the first-order Markov model, it was simplified to  $p(q_t | q_{t-1})$ 
  - However, it explains only short-term relations among notes
- Can we model the long-term relations using more complicated model?

# Toy Example

$$3 + 5 = 18$$

$$4 + 4 = 20$$

$$6 + 7 = 48$$

$$8 + 9 = 80$$

$$9 + 10 = ?$$

Note that “+” is not addition here

# Toy Example

$$3 + 5 = 18$$

$$4 + 4 = 20$$

$$y = f(x_1, x_2)$$

$$6 + 7 = 48$$

$$y = x_1 \times (x_2 + 1)$$

$$8 + 9 = 80$$

$$9 + 10 = ?$$

Note that “+” is not addition here

# Toy Example

$$2 + 2 = 6$$

$$3 + 6 = 12$$

$$4 + 5 = 19$$

$$6 + 10 = 40$$

$$7 + 18 = ?$$

Note that “+” is not addition here

# Toy Example

$$2 + 2 = 6$$

$$3 + 6 = 12 \quad y = f(x_1, x_2)$$

$$4 + 5 = 19 \quad y = \sqrt{x_1 + x_2} + {x_1}^2$$

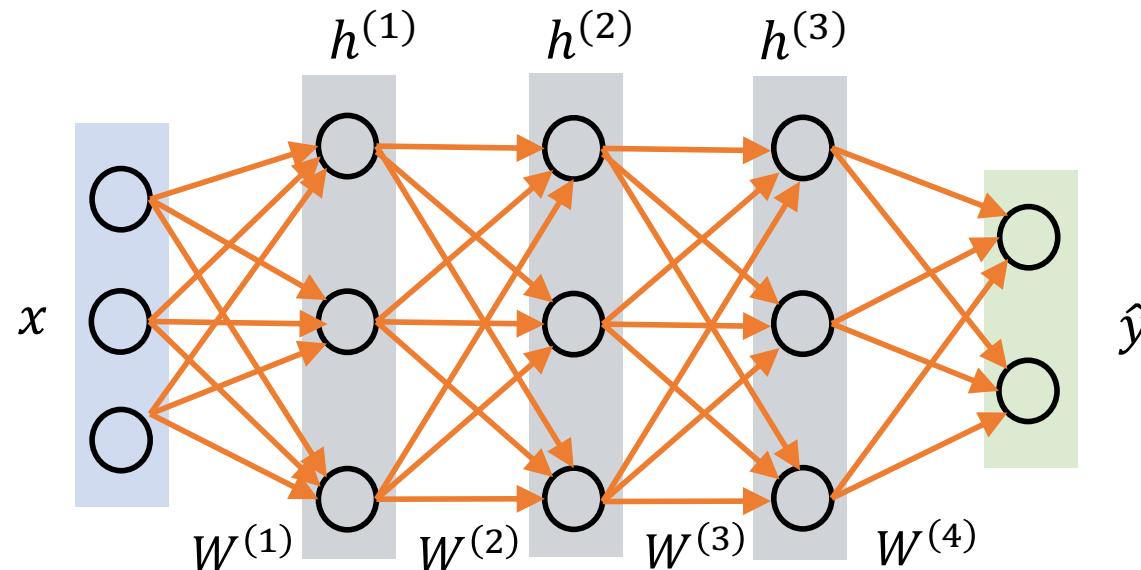
$$6 + 10 = 40$$

$$7 + 18 = ?$$

Note that “+” is not addition here

# Neural Network

- A learning model based on multi-layered networks
  - The basic model (MLP) is composed of linear transforms and element-wise nonlinear functions



Multi-Layer Perceptron (MLP)

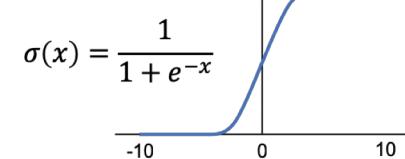
$$h^{(1)} = g^{(1)}(W^{(1)}x + b^{(1)})$$

$$h^{(2)} = g^{(2)}(W^{(2)}h^{(1)} + b^{(2)})$$

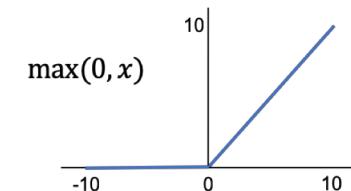
$$h^{(3)} = g^{(3)}(W^{(3)}h^{(2)} + b^{(3)})$$

$$\hat{y} = \sigma(h^{(3)})$$

Sigmoid



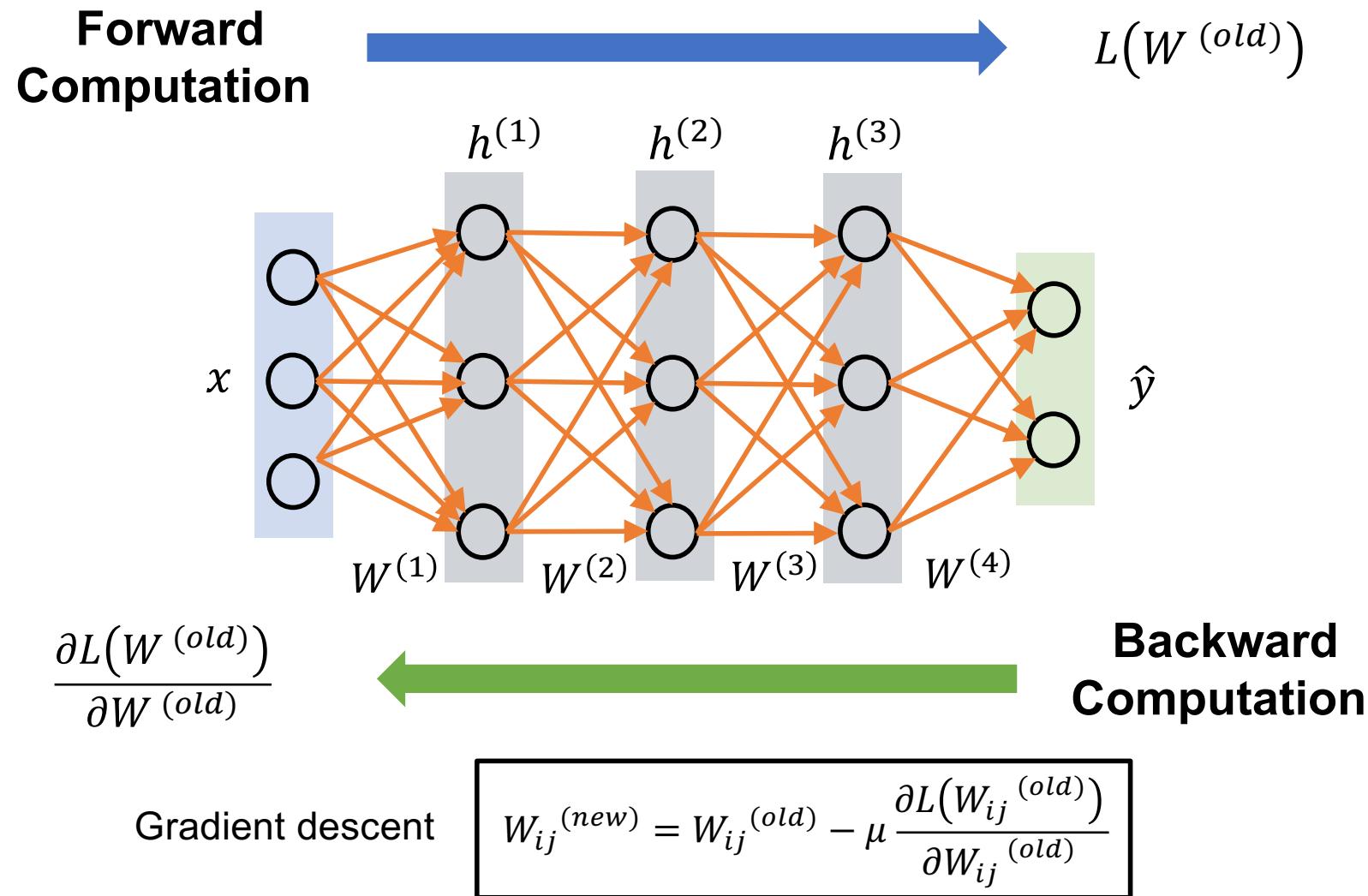
ReLU



Non-linear  
functions

# Neural Network

- The Neural network is trained via error back-propagation

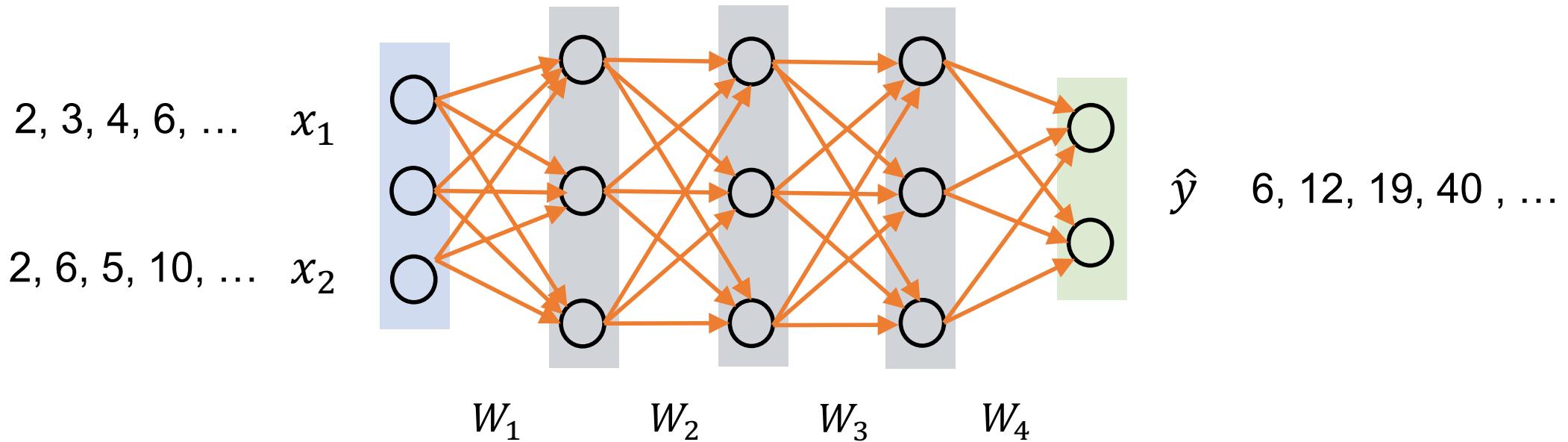


# MLP Demo and visualization

- <https://playground.tensorflow.org>

# The Toy Example

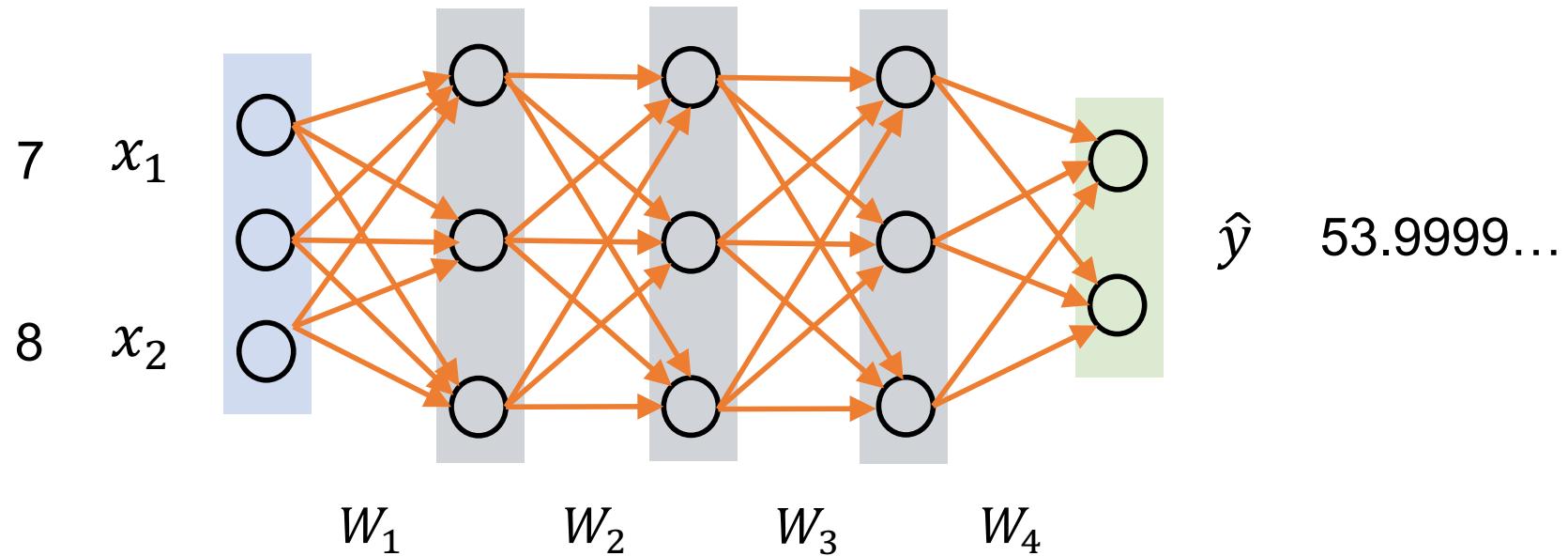
- The neural network can learn highly complicated relations between input and output



$$W_i^{new} \leftarrow W_i^{old} - \mu \frac{\partial \|\hat{y} - y\|}{\partial W_i}$$

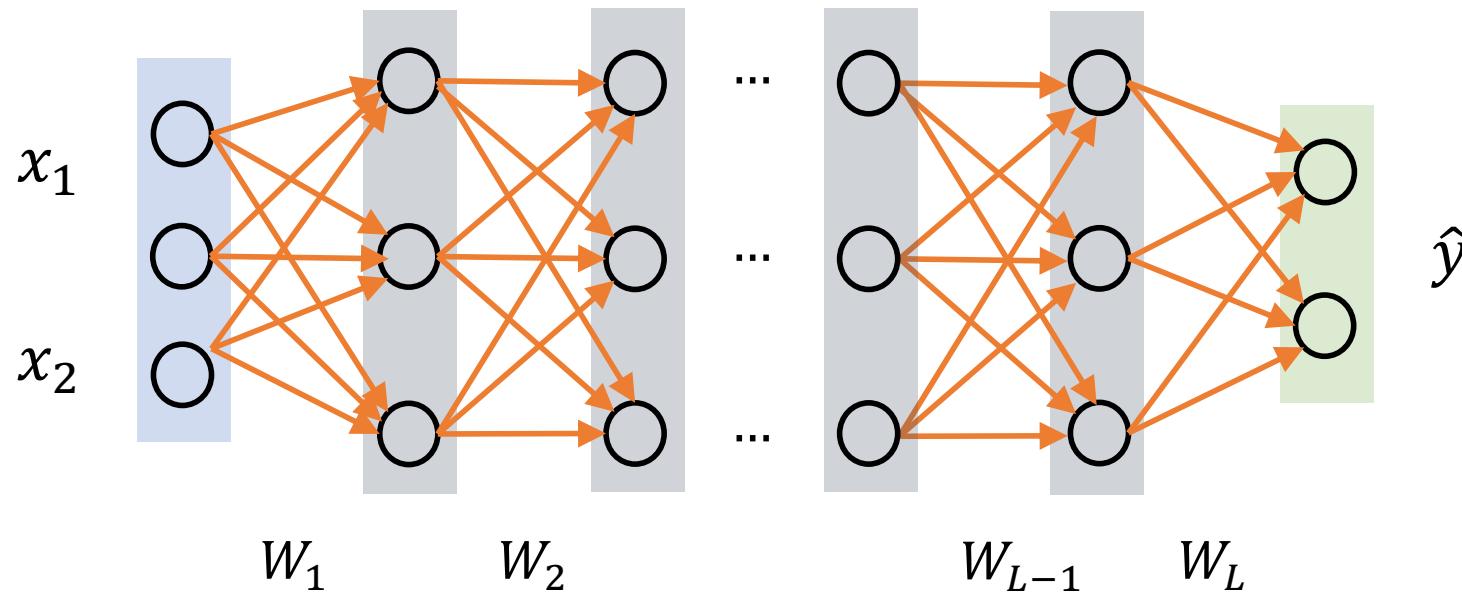
# The Toy Example

- The neural network can learn highly complicated relations between input and output



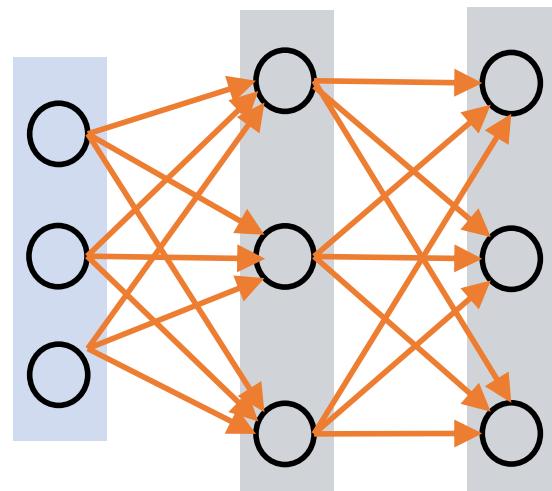
# Deep Neural Network

- Use “deep” layers
  - Many parameters to explain the data distribution
  - Need more data and fast computation (e.g. GPU)
  - Many efficient training techniques



# Deep Neural Network

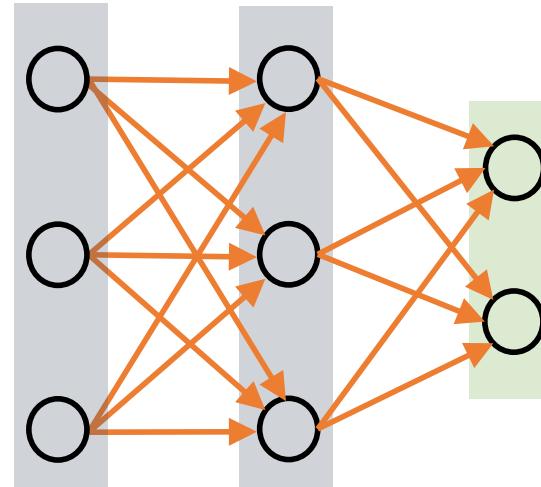
- Universal model regardless of the domain (image, audio, text, ...)

 $W_2$ 

...

...

...

 $W_{L-1}$  $W_L$ 

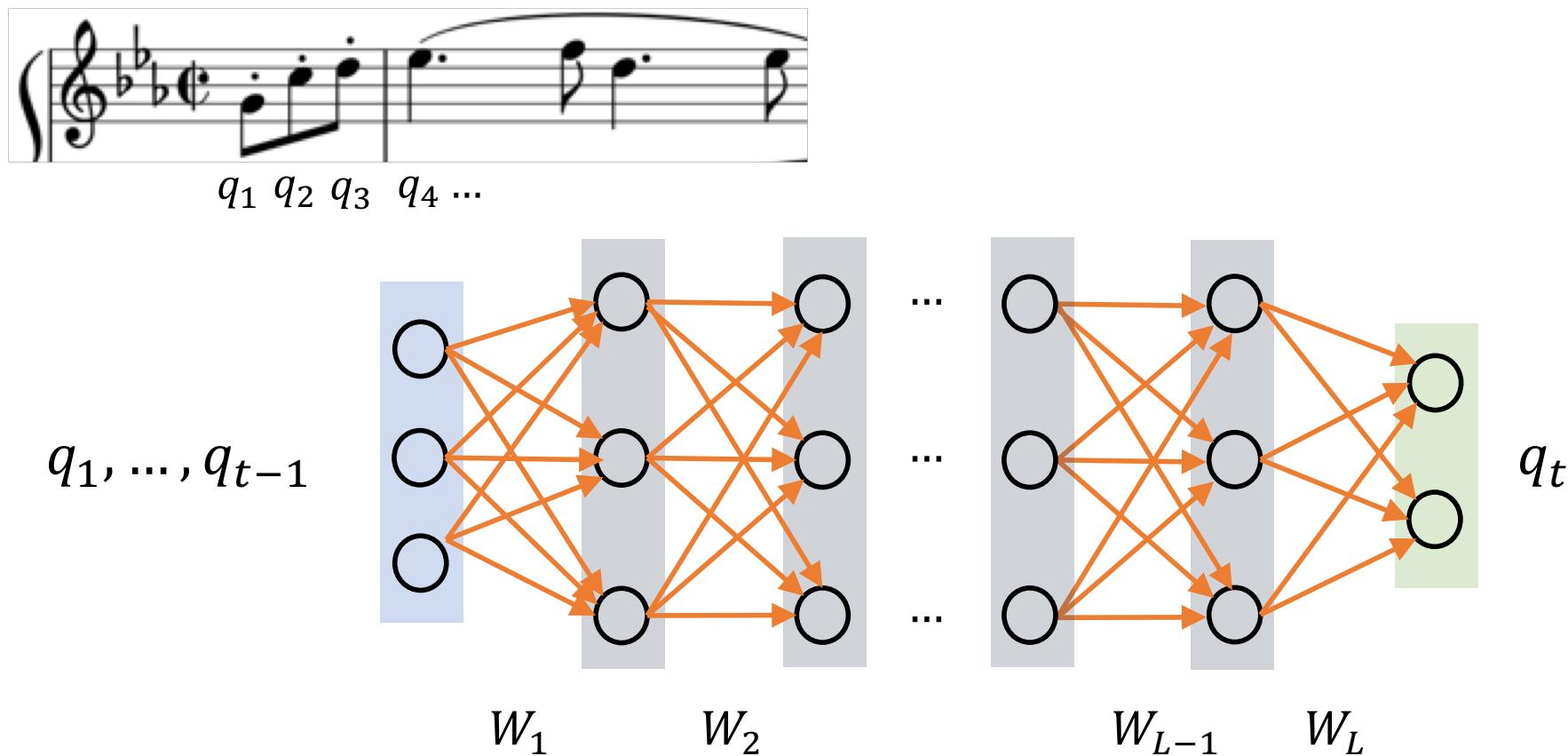
“motor-bike”

“I love coffee”

“오늘 남북정상이 만나...”

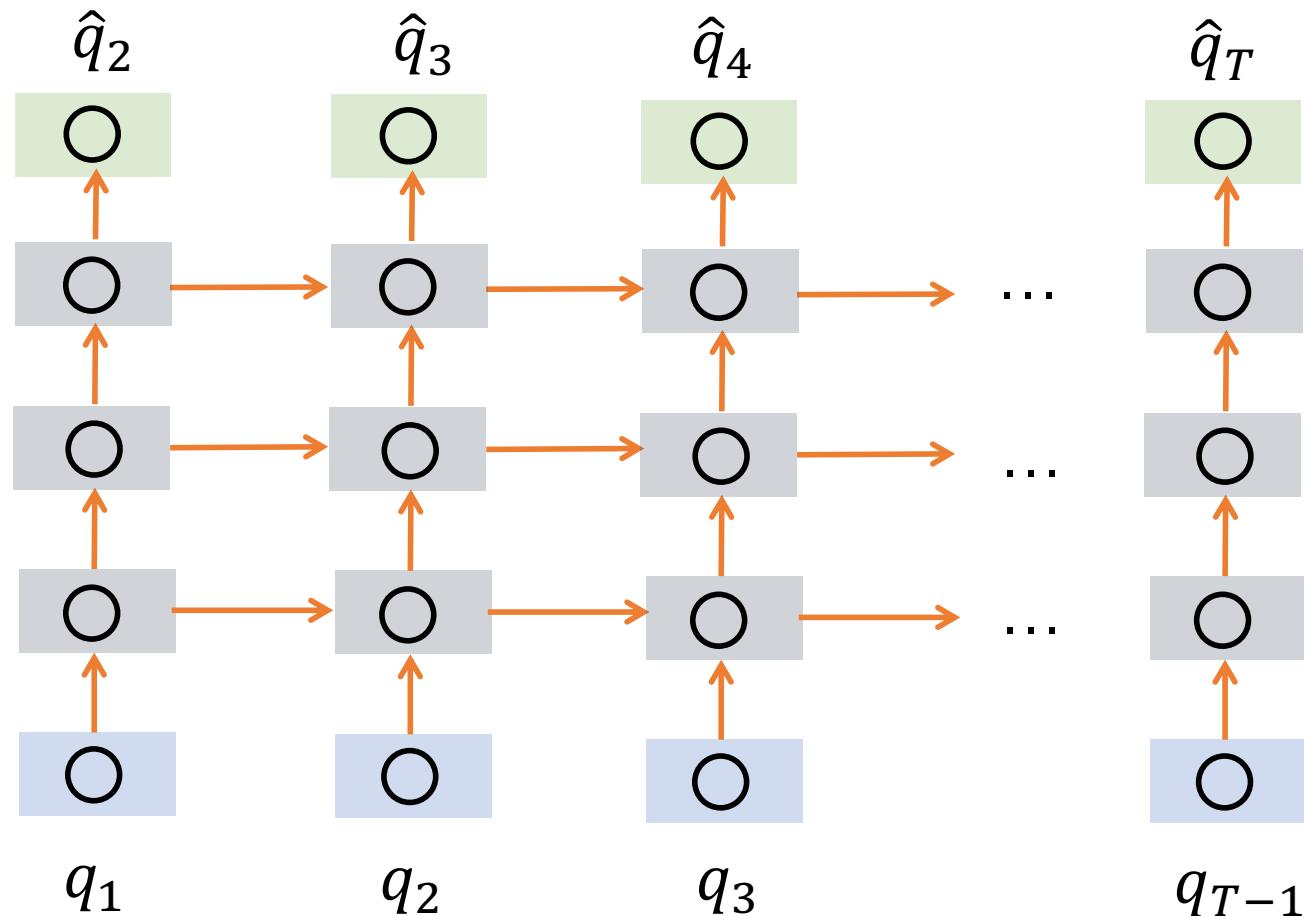
# Deep Neural Network

- Thus, we can apply the model to music!
  - However, we need to handle long sequences and variable lengths



# Recurrent Neural Networks (RNN)

- Sequence-to-sequence modeling

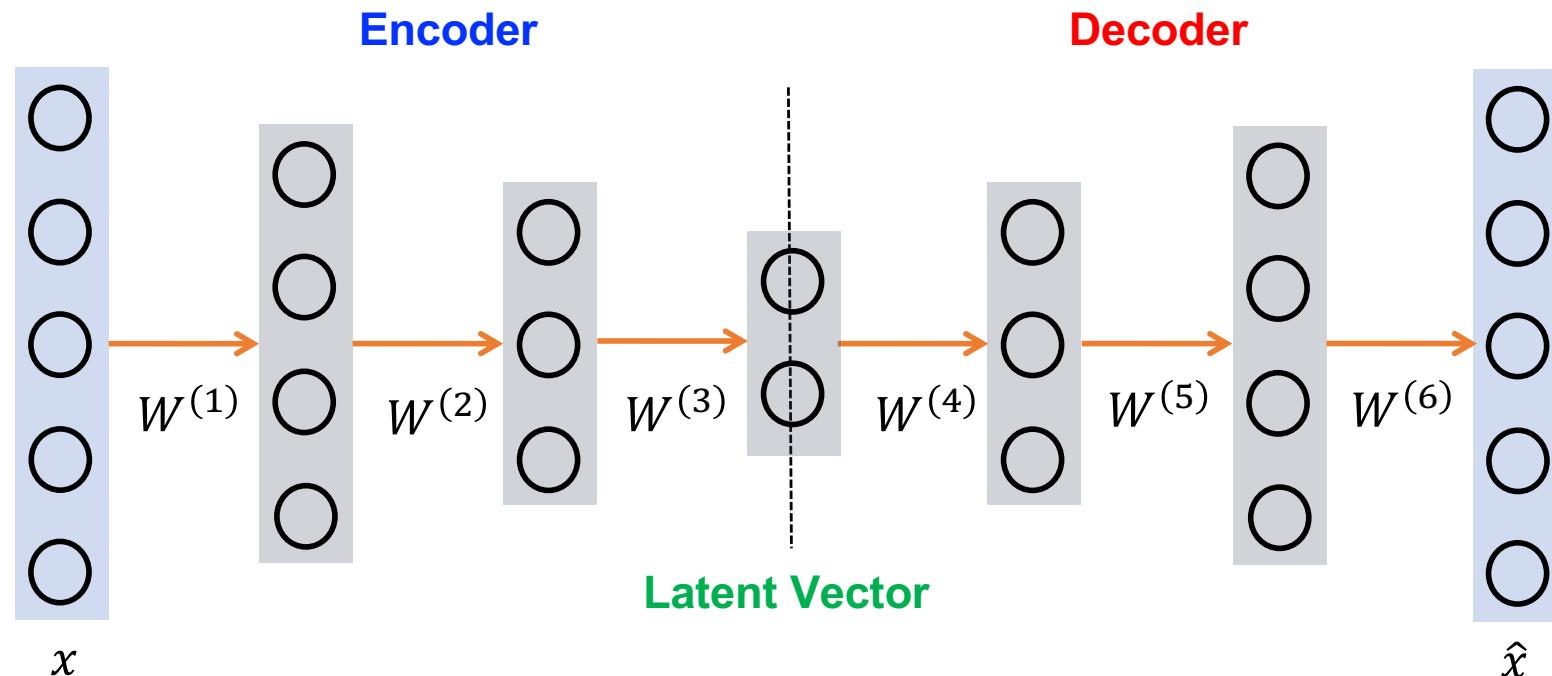


# Examples

- FolkRNN
  - <https://folkrnn.org/>
- DeepBach
  - <http://www.flow-machines.com/archives/deepbach-polyphonic-music-generation-bach-chorales/>
- DeepJazz
  - <https://deepjazz.io/>
- PerformanceRNN
  - <https://magenta.tensorflow.org/performance-rnn>

# Auto-Encoder

- Neural networks configured to reconstruct the input
  - The latent vector contains compressed information of the input
  - The decoder can be used to generate data: Variational Auto-Encoder (VAE) is more often used



Train to minimize the reconstruction error:  $L(W; x) = \|x - \hat{x}\|^2$

# Generation Examples

- Interpolation from the latent space

|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 6 | 6 | 6 | 6 | 6 | 6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 9 | 4 | 4 | 4 | 2 | 2 | 2 | 2 | 2 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 |
| 9 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 8 | 5 | 5 | 6 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 2 |
| 9 | 9 | 2 | 2 | 2 | 2 | 2 | 2 | 2 | 3 | 5 | 5 | 6 | 6 | 0 | 0 | 0 | 0 | 0 | 0 | 2 |
| 9 | 9 | 4 | 2 | 2 | 2 | 2 | 2 | 2 | 3 | 3 | 3 | 5 | 5 | 5 | 5 | 8 | 5 | 5 | 3 | 3 |
| 9 | 9 | 9 | 2 | 2 | 2 | 2 | 2 | 3 | 3 | 3 | 3 | 3 | 5 | 5 | 5 | 5 | 5 | 5 | 3 | 3 |
| 9 | 9 | 9 | 9 | 2 | 2 | 2 | 2 | 3 | 3 | 3 | 3 | 3 | 3 | 5 | 5 | 5 | 5 | 5 | 5 | 3 |
| 9 | 9 | 9 | 9 | 9 | 2 | 2 | 2 | 2 | 3 | 3 | 3 | 3 | 3 | 3 | 5 | 5 | 5 | 5 | 5 | 3 |
| 7 | 9 | 9 | 9 | 9 | 9 | 9 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 5 | 5 | 5 | 5 | 5 | 7 |
| 7 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 5 | 5 | 5 | 5 | 7 |
| 7 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 3 | 3 | 3 | 3 | 3 | 3 | 3 | 8 | 8 | 8 | 8 | 7 |
| 7 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 7 |
| 7 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 7 |
| 7 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 7 |
| 7 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 7 |
| 7 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 7 |
| 7 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 7 |
| 7 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 8 | 8 | 8 | 8 | 8 | 8 | 8 | 7 |
| 7 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 8 | 8 | 8 | 8 | 8 | 8 | 7 |
| 7 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 8 | 8 | 8 | 8 | 8 | 7 |
| 7 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 9 | 8 | 8 | 8 | 8 | 7 |
| 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 | 7 |



(Auto-Encoding Variational Bayes, Kingma and Welling, 2014)

# Music Examples

- Music-VAE
  - <https://magenta.tensorflow.org/music-vae>