# Chapter 2

# Background

fliang: Similarity to language modeling, char-rnn vs word-rnn similar so note-rnn vs chord-rnn should be as well. Benefit is that chord vocabulary MUCH larger than natural language

fliang: Try mimicking Franklin [9]

.

fliang: Motivation: Long-term dependencies are at the heart of what defines a style of music, with events spanning several notes or bars contributing to the formation of metrical and phrasal structure Cooper and Meyer [4].

fliang: Describe automatic composition vs harmonization vs analysis/scoring

Generative models are preferrable because they can be:

1. Conditioned for harmonization

2. Sampled for automatic composition

3. Queried for analysis/scoring

Automatic composition is described by Papadopoulos and Wiggins [19] as:

> a sequence (set) of rules (instructions, operations) for solving (accomplishing) a [particular] problem (task) [in a finite number of steps] of combining musical parts (things, elements) into a whole (composition)

Pearce et al. [21] emphasizes *algorithmic composition* should "be applied exclusively to activities in which the motivations have this artistic character."

fliang: Describe piano roll vs sheet music

*Piano roll* music transcriptions (**??**) are quantized both in time ($t \in T$) and note frequencies ($n \in N$). frequencies quantized to a piano roll.

fliang: Motivate quantization with Western music

.

We can represent a piano roll transcription as a high-dimensional vecctor $X_{t,n} \in \mathbb{R}^{|T| \times |N|}$ where $X_{t,n}$ denotes the note velocities for note $n$ at time $t$.

fliang: Describe MIDIs 127 iquantized pitches, isomorphic to musicXML

fliang: Describe Bach and chorales, 4 parts, soprano lead with ATB harmony, regular phrases, fermatas to denote phrase ends. Background section in chorale harmonization.

fliang: Define meter

Eck and Lapalme [7] first addressed. *Meter* is the sense of a periodic pattern of strong and weak beats which arise from periodic articulation of notes in common locations. It is implied in Western music, where bars establish perodic measures of equal length [13]. Meter provides us with key information about musical structure which can be used to predict chord changes and repepetition boundaries [4].

This chapter provides background information on two topics heavily utilized throughout this dissertation: Western music theory and recurrent neural networks. It introduce the definitions and notation used in later sections.

# 2.1   A primer on Western music theory

Music theory is a branch of musicology concerned with the study of the rules and practices of music. While the general field includes study of acoustic qualities such as dynamics and timbre, we restrict the scope of our research to modeling musical *scores* (e.g. fig. 2.1) and neglect issues related to articulation and performance (e.g. dynamics, accents, changes in tempo) as well as synthesis/generation of the physical acoustic waveforms.



Fig. 2.1 Sheet music representation of some bars from a musical score (BWV133.6) with articulation instructions removed.

This is justified because the phsyical waveforms are more closely related to the skill of the performers and instruments used and are likely to vary significantly across different performances. Furthermore, articulations in the same musical piece may differ across transcriptions and performers. Despite these variations, a piece of music is still recognizable from just the notes, suggesting that notes are the defining element for a piece of music.

### 2.1.1   Notes: the basic building blocks

A *note* is the most fundamental element of music composition and represents a sound played at a certain *pitch* for a certain *duration*. In sheet music such as fig. 2.1, the are denoted by the filled/unfilled black heads which may also possess protruding stems. As a *score* can be viewed as a collection of notes over time, notes are the fundamental building blocks for musical compositions.

#### Pitch

Though pitch is closely related to physical frequency of vibration of a waveform (as measured in Hertz), pitch its a perceptual property whose semantic meaning is derived from a listener's perception. This distinction has been scrutinized by Terhardt [27], whose visual analogy in fig. 2.2 illustrates how a pitch can be heard even if its percieved frequency is absent just as one may see the word "PITCH" despite being presented with only a suggestive shadow.
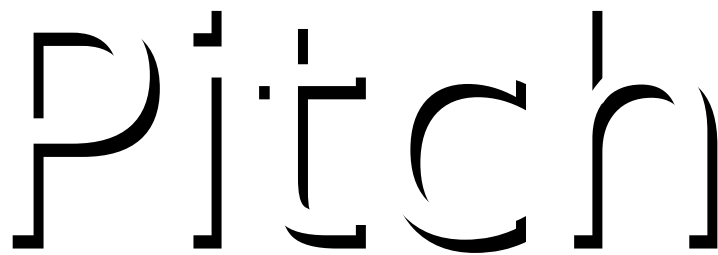


Fig. 2.2 Terhardt's visual analogy for pitch. Similar to how the viewer of this figure may percieve contours not present, pitch describes subjective information received by the listener even when physical frequencies are absent.

Despite its psychoacoustic nature, it is nevertheless useful to objectively quantify pitch as a frequency. To do so, we first need some definitions. The difference between two frequencies is called an *interval* and an *octave* is an interval corresponding to the distance between a frequency $f \in \mathbb{R}^+$ and its doubling $2f$ or halving $f/2$. Two frequencies spaced exactly an octave apart are perceived to be similar, suggesting that music is percieved on a logarithmic scale.

1   Most Western music is based on the *twelve-note chromatic scale*, which divides an *octave*
2   into twelve distinct frequencies. The *tuning system* employed dictates the precise intervals
3   between subdivisions, with *equal temperament tuning* (all subdivisions are equally spaced on
4   a logarithmic scale) presently the most widely used method [6].
5   fliang: Talk about well-tempered tuning and bach

6   Under twelve-note equal temperament tuning, the distance between two consecutive sub-
7   divisions (1/12 of an octave) is called a *semitone* (British) or *half-step* (North American) and
8   two semitones constitutes a *tone* or *whole-step*.

9   When discussing music, *note names* which enable succinct specification of a musical pitch
10  are often employed. In *scientific pitch notation*, *pitch classes* which represent a pitch modulo
11  the octave are specified by a letter ranging from *A* to *G* and optionally a single *accidental*. Pitch
12  classes without accidentals are called *natural* and correspond to the white keys on a piano. Two
13  accidentals are possible: sharps (#) raise the natural pitch class up one semitone and flats (♭)
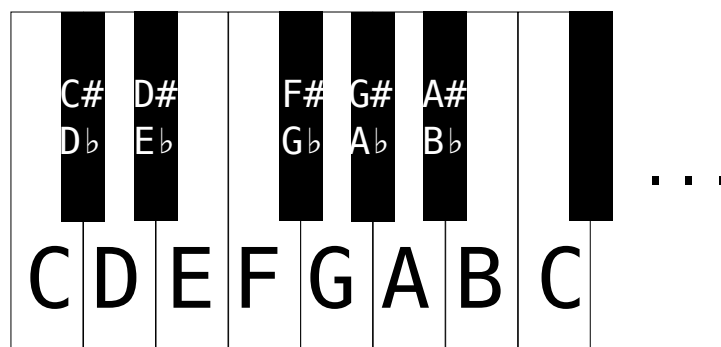14  lower by one semitone. fig. 2.3 illustrates how these pitch classes map to keys on a piano.



Fig. 2.3 Illustration of an octave in the 12-note chromatic scale on a piano keyboard.

15  Since pitch classes represent equivalence class of frequencies spaced an integral number
16  of octaves apart, unambiguously specifying a pitch requires not only a pitch class but also an
17  octave. In scientific pitch notation, this is accomplished by appending an octave number to a
18  pitch class letter (see fig. 2.4). Together, a pitch class and octave number uniquely specify the
19  notation for a pitch. On sheet music, the pitch of a note is indicated by its vertical position with
20  respect to the *stave* (the five horizontal lines and four spaces).
21  fliang: Cite wiki scientific pitch notation

## Duration

23  In addition to pitch, a note also possesses a *duration*. The duration of a note indicates how long
24  it is to be played and is measured in fractions of a *whole note* (American) or *semibreve* (British).
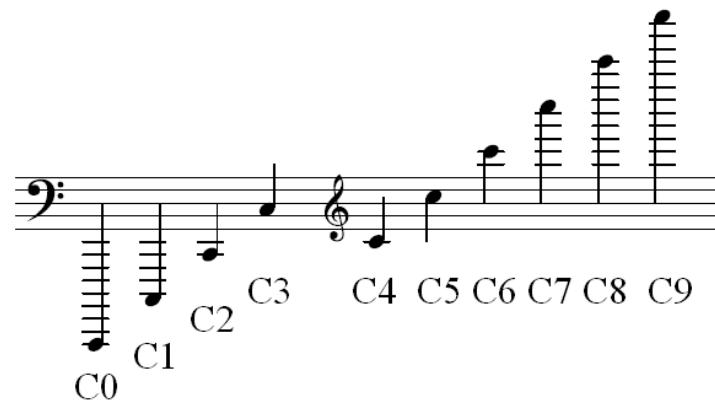
Fig. 2.4 Scientific pitch notation and sheet music notation of *C* notes at ten different octaves.

Perhaps the most common duration is a *quarter-note* (American) or *crotchet* (British). Other note durations are also possible and the most common along with their notation in sheet music are enumerated in fig. 2.5. The relationship between durations and physical time intervals is given by the *tempo*, which is usually denoted near the start of the piece in beats per minute.

fliang: Cite wiki Whole note



Fig. 2.5 Comparison of various note durations.

## 2.1.2   Key signature

*Tonal music* is a genre of music encompassing most of Western classical music. It is characterized by the prevalence of one pitch class (the *tonic*) around which the remainder of the piece is built.

A basic concept within tonal music is *mode*, which defines a subset of pitch classes that are "in key" with respect to the tonic. table 2.1 illustrates the pitch intervals between adjacent pitch classes within two important modes: the *major* (Ionian, I) and *minor* (Aeolian, VI) modes. The choice of tonic and mode is collectively referred to as the *key signature*.

| Mode | Pitch Intervals (semitones) |
|------|------------------------------|
| Major (Ionian, I) | +2, +2, +1, +2, +2, +2 |
| Minor (Aeolian, VI) | +2, +1, +2, +2, +1, +2 |

Table 2.1 Pitch intervals for the two most important modes[10]. The pitches in a scale can be found by starting at the tonic and successively offsetting by the given pitch intervals.

### 2.1.3   Polyphony, chords, and chord progressions

Whereas *monophonic* music is characterized by the presence of a single *part* sounding at most one note at any given time, *polyphonic* music contains multiple parts potentially sounding multiple pitches at the same time. Just as notes form the basis of monophonic music, chords are the fundamental building blocks for polyphonic music.

**Chords: basic units for representing simultaneously sounding notes**

A *chord* is a collection of three or more pitches all sounding simultaneously[23]. In Western classical music, they typically consist of a *root note* whose pitch class forms a base from which successive notes are built upon. The intervals between the pitch classes in a chord are commonly labeled using *qualities*, which are invariant across octaves. Different realizations of the same chord (e.g. octave choices for each pitch class) are called *voicings*.

table 2.2 lists some common chord qualities and their corresponding intervals from the root note. Chord names are given as a root pitch class followed by a quality, for example: *C* major, *A* minor, or *G* half-diminished 7th.

The lowest note in a chord is called the *bass* note and is oftentimes the root note. However, alternative voicings called *inversions* can place the root note on a different octave and cause the bass and root notes to differ.

**Chord progressions and phrases**

Sequences of chords are called *chord progressions*, which are oftentimes grouped with adjacent progressions within a score into coherent units called *phrases*. Many psychoacoustic phenomena such as stability, mood, and expectation can be attributed to phrase structure and

| Chord quality | Intervals from root pitch class |
|---|:---:|
| Major | +4, +7 |
| Major 6th | +4, +7, +8 |
| Major 7th | +4, +7, +11 |
| Minor | +3, +7 |
| Minor 6th | +3, +7, +9 |
| Minor 7th | +3, +7, +10 |
| Dominant 7th | +4, +7, +10 |
| Augmented | +4, +8 |
| Diminished | +3, +6 |
| Diminished 7th | +3, +6, +9 |
| Half-diminished 7th | +3, +6, +10 |

Table 2.2 Common chord qualities and their corresponding intervals[10]

choice of chord progressions. For example, chord progressions called *harmonic cadences* are    1
commonly used to conclude a phrase and create a sense of resolution[23]. Another impor-    2
tant example are *modulations*, where a chord progression is used to transition the music into a    3
different key signature.    4

    As chords can be overlapping and contain notes straddling between two chords or involve    5
uncommon chord qualities, identifying chord progresssions involves a degree of subjectivity. A    6
common method for analyzing chord progressions is *roman numeral analysis*, where I denotes    7
the tonic pitch class,successive roman numerals correspond to successive pitch classes in the    8
key, and capitalization is used to denote major/minor. For example, the chord progression $C$    9
major – $A$ minor – $D$ major 7th – $G$ major in the $C$ major key signature would be represented    10
as $I - ii - IImaj7 - V$.    11

**Transposition invariance**    12

A common theme throughout our discussion has been ambiguity of the tonic. When discussing    13
key signatures, the mode was defined using intervals relative to a choice of tonic. Similarly,    14
roman numeral analysis of chord progressions is also conducted relative to a tonic. Even if a    15
chord progression and tonic are both transposed by arbitrary pitch interval, the roman numeral    16
analysis will remain unchanged.    17

    The *transposition invariance* of chord progressions and modes is an important property of    18
music. It enables us to offset an entire score by an arbitrary pitch interval without affecting the    19
important psychoacoustic qualities captured by chord progressions and choice of mode.    20

## ₁ 2.2    Neural sequence probability modeling

₂ Our work in later sections make heavy use of neural networks. In this section, we briefly review
₃ the relevant concepts and set up notation.

### ₄ 2.2.1    Neurons: the basic computation unit

₅ Neurons are the basic abstraction which are combined together to form neural networks. A
₆ *neuron* is a parametric model of a function $f : \mathbb{R}^D \to \mathbb{R}$ from its $D$-dimensional input $\boldsymbol{x}$ to its
₇ output $y$. Our neurons will be defined as

$$f(\boldsymbol{x}) \coloneqq \sigma(\langle \boldsymbol{w}, \boldsymbol{x} \rangle) \tag{2.1}$$

₉ which can be viewed as an inner product with *weights* $\boldsymbol{w}$ to produce an *activation* $z \coloneqq \langle \boldsymbol{w}, \boldsymbol{x} \rangle \in$
₁₀ $\mathbb{R}$ which is then squashed to a bounded domain by a non-linear **activation function** $\sigma : \mathbb{R} \to$
₁₁ $[L, U]$. This is visually depicted in fig. 2.6, which also makes apparent the interpretation of
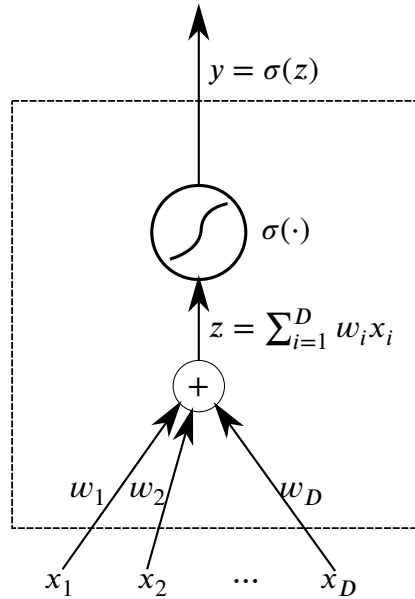₁₂ weight $w_i$ as the sensitivity of the output $y$ to the input $x_i$.



Fig. 2.6 A single neuron first computes an activation $z$ and then passes it through an activation function $\sigma(\cdot)$

### 2.2.2   Feedforward neural networks

Multiple neurons may share inputs and have their outputs concatenated together to form a *layer*
modelling a multivariate functions $f : \mathbb{R}^{D_{\text{in}}} \to \mathbb{R}^{D_{\text{out}}}$. Multiple layers can then be composed
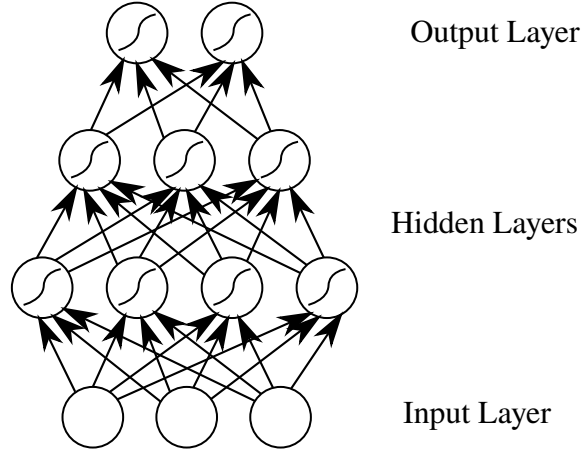together to form a *feedforwd neural network*.



Fig. 2.7 Graph depiction of a feedforward neural network with 2 hidden layers

Although a single hidden layer is theoretically sufficient for a universal function approximator[5],
the number of hidden units to guarantee reported theoretical bounds are usually unfeasibly
large. Instead, recent work in *deep learning* has shown that deep models which contain many
hidden layers can achieve strong performance across a variety of tasks[1].

The improved modeling capacity gained by composing multiple layers is due to the com-
position of multiple non-linear activation functions. In fact, it is easy to show that removing
activation functions would make a deep network equivalent to a single matrix transform: let
$\boldsymbol{W}_{l,l+1}$ denote the weights between layers $l$ and $l + 1$. The original neural network computes
the function

$$\sigma \left( \boldsymbol{W}_{L,L-1} \sigma \left( \boldsymbol{W}_{L-1,L-2} \cdots \sigma \left( \boldsymbol{W}_{2,1} \boldsymbol{x} \right) \cdots \right) \right) \tag{2.2}$$

After removing the activation functions $\sigma$, we are left with

$$\boldsymbol{W}_{L,L-1} \boldsymbol{W}_{L-1,L-2} \cdots \boldsymbol{W}_{2,1} \boldsymbol{x} = \boldsymbol{x} = \tilde{\boldsymbol{W}} \boldsymbol{x} \tag{2.3}$$

where $\tilde{\boldsymbol{W}} = \left( \prod_{i=1}^{L-1} \boldsymbol{W}_{i,i+1} \right)$ is a matrix transform computing the same function as the neural
network with activation functions removed.

### 2.2.3   Recurrent neural networks

While feedforward neural networks provide a flexible model for approximating arbitrary functions, they require a fixed-dimension input $x$ and hence cannot be directly applied to sequential data $x = (x_t)_{t=1}^{T}$ where $T$ may vary.

A naive method for extending feedforward networks would be to independently apply a feedforward network to compute $y_t = f(x_t \theta)$ at each timestep $1 \leq t \leq T$. However, this approach is only correct when each output $y_t$ depends only on the input at the current time $x_t$ and is independent of all prior inputs $\{x_k\}_{k<t}$. This assumption is false in musical data: the current musical note usually is highly dependent on the sequence of notes leading up to it.

This shortcoming motivates *recurrent neural networks* (RNNs), which generalize feedforward networks by introducing time-delayed recurrent connections between hidden layers (Elman networks [8]) or from the output layers to the hidden layers (Jordan networks [16]). Mathematically, an (Elman-type) RNN is a discrete time dynamical system commonly parameterized as:

$$h_t = W_{xh}\sigma_{xh}(x_t) + W_{hh}\sigma_{hh}(h_{t-1}) \tag{2.4}$$

$$y_t = W_{hy}\sigma_{hy}(h_t) \tag{2.5}$$

where $\sigma_{..}(\cdot)$ are activation functions acting element-wise and $\theta = \{W_{xh}, W_{hh}, W_{hy}\}$ are the learned parameters. fig. 2.8 provides a graphical illustration of such a network. Notice that apart from the edges between hidden nodes, the network is identical to a regular feedforward network (fig. 2.7).
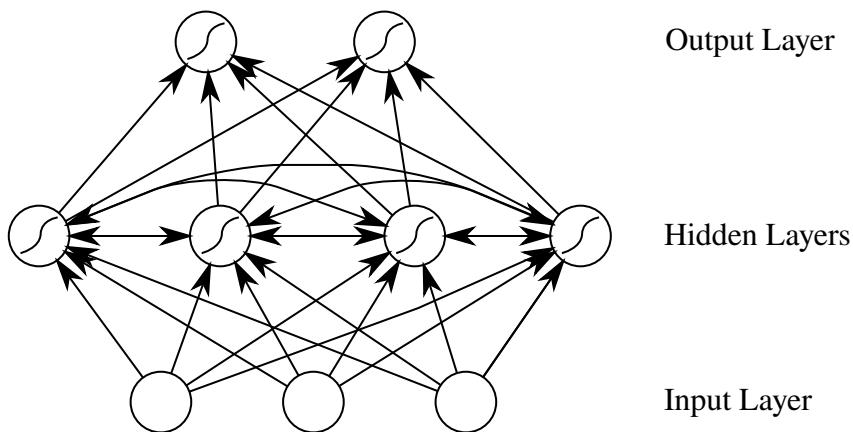
> fliang: Why do we use Elman



Fig. 2.8 Graph representation of an Elman-type RNN.

To apply the RNN over an input sequence $\boldsymbol{x}$, the activations of the hidden states are first    1
initialized to an initial value $\boldsymbol{h} \in \mathbb{R}^{D_h}$. Next, for each timestep $t$ the hidden layer activations    2
are computed using the current input $\boldsymbol{x}_t$ and the previous hidden state activations $\boldsymbol{h}_{t-1}$. This    3
motivates an alternative perspective on RNNs as a template consisting of a feedforward network    4
with inputs $\{\boldsymbol{x}_t, \boldsymbol{h}_{t-1}\}$ (see fig. 2.9) replicated across time $t$.    5
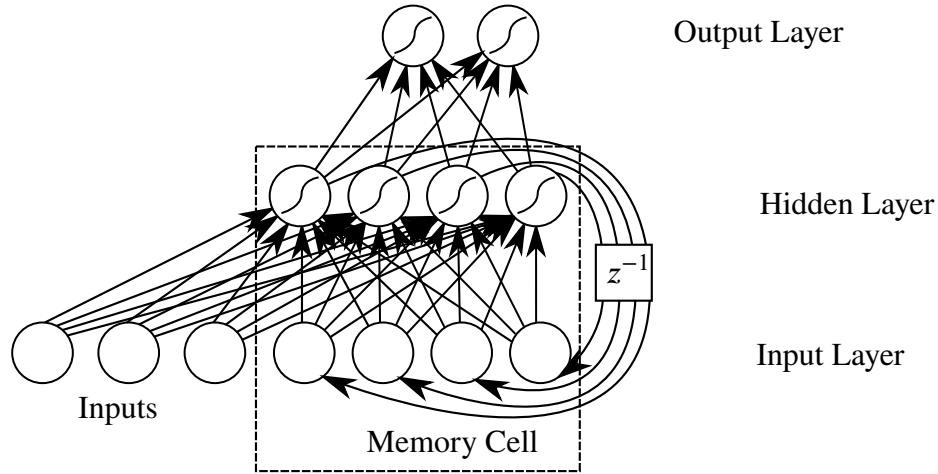


Fig. 2.9 Equivalent formulation of an Elman-type RNN treating the time-delayed hidden state $\boldsymbol{h}_{t-1}$ as additional inputs to a feedforward network

**Memory cell abstraction**    6

Notice that fig. 2.9 delineates the recurrent hidden state from the rest of the diagram, intro-    7
ducing an abstraction called a *memory cell*. This allows us to abstract away how $\boldsymbol{y}_t$ and $\boldsymbol{h}_t$    8
are computed from $\boldsymbol{x}_t$ and $\boldsymbol{h}_{t-1}$, enabling discussion of RNNs applicable to many different    9
implementations. Concretely, A memory cell is an interface which for each timestep $t$:    10

- Takes inputs    11

    - The current element in the input sequence $\boldsymbol{x}_t$    12

    - The previous hidden state $\boldsymbol{h}_{t-1}$    13

- Produces outputs    14

    - The next hidden state $\boldsymbol{h}_t = f_h(\boldsymbol{x}_t, \boldsymbol{h}_{t-1})$    15

    - The outputs $\boldsymbol{y}_t = f_y(\boldsymbol{h}_t)$.    16

In future diagrams, we will abstractly represent the memory cell as a node labelled with $\boldsymbol{h}$.    17

₁ **Unrolling into a DAG**

₂ One important operation on RNNs is *unrolling* the template in fig. 2.9 into a chain of $T$ repli-
₃ cations with connected hidden states (fig. 2.10). This removes the time-delayed recurrent,
₄ converting the RNN into a finite directed acyclic graph where nodes represent pieces of data
₅ and edges $s \rightarrow t$ indicate that $t$ is a function of $s$. This is identical to a feedforward network,
₆ justifying the view of RNNs as a dynamically-sized feedforward network with $T$ layers and
₇ parameters tied across all layers.

₈ > fliang: Talk about how view of unrolled RNN as DNN inspires cross-polination e.g. gating to let information flow deeper: LSTM gates and highway networks, grid LSTMS
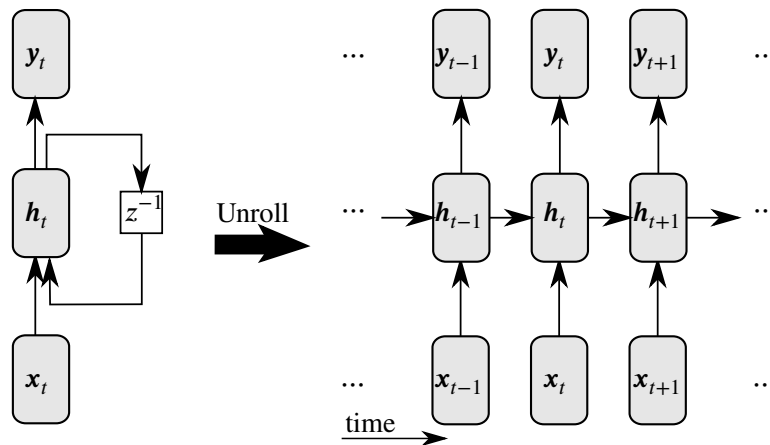


Fig. 2.10 Signal flow diagram representation of a single-layer RNN and its corresponding directed acyclic graph after unrolling

₉      fig. 2.10 makes it obvious how the hidden state is carried along throughout the sequence of
₁₀ computations, giving rise to a useful alternative interpretation of the hidden state as a temporal
₁₁ memory mechanism. Under this interpretation, we can view the hidden state update $h_t = $
₁₂ $f_h(x_t, h_{t-1})$ as *writing* information extracted from the current inputs $x_t$ to the memory $h_{t-1}$.
₁₃ Similarly, producing the outputs $y_t = f_y(h_t)$ can be seen as *reading* information from the
₁₄ hidden state.

₁₅ > fliang: Compare to N-grams; show how it's like an infinite context. One interpretation is to view the hidden state $h_t$ as an infinite-length prior context window, summarizing all of the prior inputs into into a compact fixed-size vector.

**Stacking memory cells to form multi-layer RNNs**                          1

Since the RNN outputs $y$ also form a sequence with the same length as the inputs $x$, they can   2
be used as inputs into another RNN. This stacking of multiple memory cells is similar to the   3
layering seen in deep neural networks, giving rise to the term *deep neural sequence models*   4

fliang: Cite?                                                                5

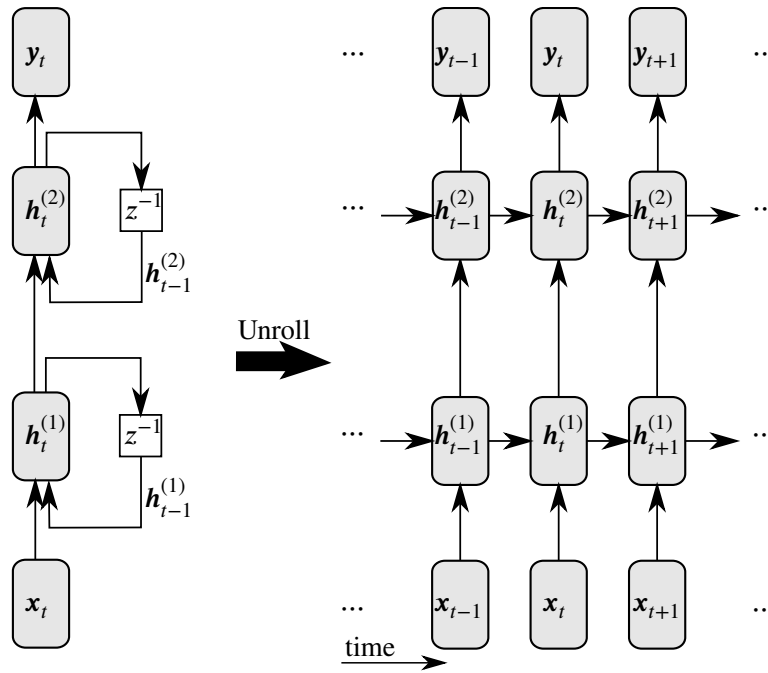. This is illustrated in fig. 2.11.                                          6



Fig. 2.11 Template and unrolling of a stacked 2-layer RNN

The greater modeling capabilities of multilayer RNNs can be attributed to two primary fac-   7
tors: composition of multiple non-linearities and an increase in the number of paths through   8
which information can flow. The former is analogous to the feedforward case: stacking mem-   9
ory cells increases the number of non-linearities in the composite cell just like stacking multiple   10
layers in feedforward networks. To understand the latter point, notice that in fig. 2.10 there is   11
only a single path from $x_{t-1}$ to $y_t$ hence the conditional independence $y_t \perp\!\!\!\perp x_{t-1} | h_t^{(1)}$ is satis-   12
fied. However, in fig. 2.11 there are multiple paths from $x_{t-1}$ to $y_t$ (e.g. passing through either   13
$h_{t-1}^{(2)} \rightarrow h_t^{(2)}$ or $h_{t-1}^{(1)} \rightarrow h_t^{(1)}$) through which information may flow. Additionally, the hidden   14
state transitions occur on two seperate memory cells so parameters need not be tied and the   15
stacked RNN can learn different time dynamics at each depth.                  16

## 2.2.4 Modeling assumptions

fliang: Integrate better

We make the following assumptions about the sequence $\boldsymbol{x}_{1:T}$, $\boldsymbol{y}_{1:T}$, and $\boldsymbol{h}_{0:T}$:

1. Modified Markov assumption:

$$\forall t \; : \; P(\boldsymbol{h}_t|\boldsymbol{h}_{0:t-1}, \boldsymbol{x}_{1:t}) = P(\boldsymbol{h}_t|\boldsymbol{h}_{t-1}, \boldsymbol{x}_t) \tag{2.6}$$

2. Hidden State Stationarity:

$$\forall t_1, t_2 \; : \; P(\boldsymbol{h}_{t_1} = \boldsymbol{k}|\boldsymbol{h}_{t_1-1} = \boldsymbol{i}, \boldsymbol{x}_{t_1} = \boldsymbol{j}) = P(\boldsymbol{h}_{t_2} = \boldsymbol{k}|\boldsymbol{h}_{t_2-1} = \boldsymbol{i}, \boldsymbol{x}_{t_2} = \boldsymbol{j}) \tag{2.7}$$

3. Output Stationarity:

$$\forall t_1, t_2 \; : \; P(\boldsymbol{y}_{t_1} = \boldsymbol{j}|\boldsymbol{h}_{t_1} = \boldsymbol{i}) = P(\boldsymbol{y}_{t_2} = \boldsymbol{j}|\boldsymbol{h}_{t_2} = \boldsymbol{i}) \tag{2.8}$$

4. Output independence:

$$P(\boldsymbol{y}_{1:T}|\boldsymbol{h}_{0:T}, \boldsymbol{x}_{1:T}) = \prod_{t=1}^{T} P(\boldsymbol{y}_t|\boldsymbol{h}_t, \boldsymbol{x}_t) \tag{2.9}$$

These assumptions imply the sequential factorization:

$$P(\boldsymbol{y}_{1:T}, \boldsymbol{h}_{1:T}|\boldsymbol{h}_0, \boldsymbol{x}_{1:T}) \tag{2.10}$$

$$= P(\boldsymbol{y}_{1:T}|\boldsymbol{h}_{0:T}, \boldsymbol{x}_{1:T})P(\boldsymbol{h}_{1:T}|\boldsymbol{h}_0, \boldsymbol{x}_{1:T}) \tag{2.11}$$

$$= \left(\prod_{t=1}^{T} P(\boldsymbol{y}_t|\boldsymbol{h}_t)\right) P(\boldsymbol{h}_{1:T}|\boldsymbol{h}_0, \boldsymbol{x}_{1:T}) \qquad \text{eq. (2.9)} \tag{2.12}$$

$$= \left(\prod_{t=1}^{T} P(\boldsymbol{y}_t|\boldsymbol{h}_t)\right)\left(\prod_{t=1}^{T} P(\boldsymbol{h}_t|\boldsymbol{h}_{0:t-1}, \boldsymbol{x}_{1:t})\right) \tag{2.13}$$

$$= \left(\prod_{t=1}^{T} P(\boldsymbol{y}_t|\boldsymbol{h}_t)\right)\left(\prod_{t=1}^{T} P(\boldsymbol{h}_t|\boldsymbol{h}_{t-1}, \boldsymbol{x}_t)\right) \qquad \text{eq. (2.6)} \tag{2.14}$$

$$= \prod_{t=1}^{T} P(\boldsymbol{y}_t|\boldsymbol{h}_t, \boldsymbol{x}_t)P(\boldsymbol{h}_t|\boldsymbol{h}_{t-1}, \boldsymbol{x}_t) \tag{2.15}$$

$$\tag{2.16}$$

eq. (2.7) combined with eq. (2.8) imply that $P(\boldsymbol{y}_t|\boldsymbol{h}_t)$ and $P(\boldsymbol{h}_t|\boldsymbol{h}_{t-1}, \boldsymbol{x}_t)$ are time-invariant and
can be modelled by the same recurrent function.

In RNNs, the hidden state dynamics $P(\boldsymbol{h}_t|\boldsymbol{h}_{t-1}, \boldsymbol{x}_t)$ are deterministic:

$$\boldsymbol{h}_t = f_h(\boldsymbol{x}_t, \boldsymbol{h}_{t-1}) \tag{2.17}$$

Which means that $P(\boldsymbol{y}_{1:T}, \boldsymbol{h}_{1:T}|\boldsymbol{h}_0, \boldsymbol{x}_{1:T}) = P(\boldsymbol{y}_{1:T}|\boldsymbol{h}_0, \boldsymbol{x}_{1:T})$. This yields the factorization

$$P(\boldsymbol{y}_{1:T}|\boldsymbol{h}_0, \boldsymbol{x}_{1:T}) = P(\boldsymbol{y}_{1:T}, \boldsymbol{h}_{1:T}|\boldsymbol{h}_0, \boldsymbol{x}_{1:T}) = \prod_{t=1}^{T} P(\boldsymbol{y}_t|\boldsymbol{h}_t, \boldsymbol{x}_t) f_h(\boldsymbol{x}_t, \boldsymbol{h}_{t-1}) \tag{2.18}$$

> fliang: Draw PGM

However, one minor problem remains. Let $\boldsymbol{z}_t = f_y(f_h(\boldsymbol{x}_t, \boldsymbol{h}_{t-1}))$ (with $f_y$ and $f_h$ as defined
in

> fliang: ref

) denote the outputs of the RNN model at time $t$. Note that $\boldsymbol{z}_t$ can be any real vector in $\mathbb{R}^{|V|}$

> fliang: Define $V$ to be the vocabulary

, but $P(\boldsymbol{x}_{t+1}|\boldsymbol{h}_{t-1}, \boldsymbol{x}_t)$ is a probability vector constrained to sum to one.

Fortunately, we can treat $\boldsymbol{z}_t$ as the *scores* for a *Boltzmann distribution* (aka softmax):

$$P(\boldsymbol{y}_t = s|\boldsymbol{h}_{t-1}, \boldsymbol{x}_t) = \frac{\exp\left(-\boldsymbol{z}_{t,s}/T\right)}{\sum_{k=1}^{K}\left(\exp -\boldsymbol{z}_{t,k}/T\right)} \tag{2.19}$$

where $T \in \mathbb{R}^+$ is a *temperature* parameter (set to $T = 1$ during training and varied during sampling). To keep notation compact, we omit writing this explicitly and understand $P(\boldsymbol{y}_t|\boldsymbol{h}_{t-1}, \boldsymbol{x}_t)$
to mean the Boltzmann distribution parameterized by the scores $f_y(f_h(\boldsymbol{x}_t, \boldsymbol{h}_{t-1}))$.

Note the similarity between eq. (2.6)–eq. (2.9) and the assumptions for Hidden Markov
models [22]. Discrepancies are due to the presence of an input sequence $\boldsymbol{x}_{1:T}$ in our sequence-
to-sequence model.

> fliang: Discuss validity of assumptions, namely output independence assuming hidden state and
> input summarize all prior context

## 2.2.5   Training RNNs: backpropogation through time

The parameters $\boldsymbol{\theta}$ of a RNN are typically learned from data to minimize a cost $\mathscr{E} = \sum_{1 \leq t \leq T} \mathscr{E}_t(\boldsymbol{x}_t)$ measuring the performance of the network on some task. This optimization is commonly carried out using iterative gradient descent methods, which require computation of the gradients $\frac{\partial \mathscr{E}}{\partial \boldsymbol{\theta}}$ at each iteration.

One approach for computing the necessary gradients is *backpropagation through time* (BPTT)[12], an adaptation of the backpropagation algorithm[17, 24] to the unrolled RNN computation graph. We can apply the chain rule to the unrolled RNN's computation graph in fig. 2.12 to obtain

$$\frac{\partial \mathscr{E}}{\partial \boldsymbol{\theta}} = \sum_{1 \leq t \leq T} \frac{\partial \mathscr{E}_t}{\partial \boldsymbol{\theta}} \tag{2.20}$$

$$\frac{\partial \mathscr{E}_t}{\partial \boldsymbol{\theta}} = \sum_{1 \leq k \leq t} \left( \frac{\partial \mathscr{E}_t}{\partial \boldsymbol{y}_t} \frac{\partial \boldsymbol{y}_t}{\partial \boldsymbol{h}_t} \frac{\partial \boldsymbol{h}_t}{\partial \boldsymbol{h}_k} \frac{\partial \boldsymbol{h}_k}{\partial \boldsymbol{\theta}} \right) \tag{2.21}$$

$$\frac{\partial \boldsymbol{h}_t}{\partial \boldsymbol{h}_k} = \prod_{t \geq i > k} \frac{\partial \boldsymbol{h}_i}{\partial \boldsymbol{h}_{i-1}} = \prod_{t \geq i > k} \boldsymbol{W}_{hh}^{\mathsf{T}} \operatorname{diag}\left( \sigma'_{hh}(\boldsymbol{h}_{i-1}) \right) \tag{2.22}$$
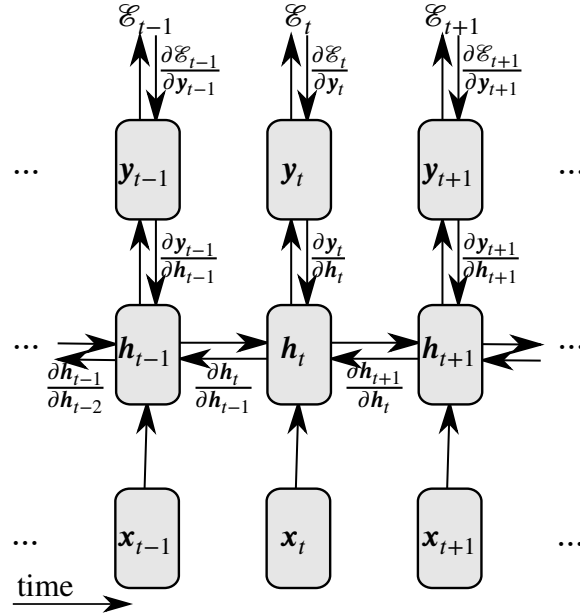


Fig. 2.12 The gradients passed along network edges during BPTT.

section 2.2.5 expresses how the error $\mathscr{E}_t$ at time $t$ is a sum of *temporal contributions* $\frac{\partial \mathscr{E}_t}{\partial \boldsymbol{y}_t} \frac{\partial \boldsymbol{y}_t}{\partial \boldsymbol{h}_t} \frac{\partial \boldsymbol{h}_t}{\partial \boldsymbol{h}_k} \frac{\partial \boldsymbol{h}_k}{\partial \boldsymbol{\theta}}$ measuring how $\boldsymbol{\theta}$'s impact on $\boldsymbol{h}_k$ affects the cost at some future time $t > k$. The factors in section 2.2.5 measures the affect of the hidden state $\boldsymbol{h}_k$ on some future state $\boldsymbol{h}_t$

where $t > k$ and can be interpreted as transferring the error "in time" from step $t$ back to step $k$ [20].

### Vanishing/exploding gradients

Unfortunately, naive implementations of section 2.2.3 and section 2.2.3 oftentimes suffers from two well known problems: the *vanishing gradient* and *exploding gradient*[2]. Broadly speaking, these problems are both related to the product in section 2.2.5 exponentially growing or shrinking for long timespans (i.e. $t \gg k$).

Following Pascanu et al. [20], let $\|\cdot\|$ be any submultiplicative matrix norm (e.g. Frobenius, spectral, nuclear, Shatten $p$-norms). Without loss of generality, we will use the *operator norm* defined as

$$\|A\| = \sup_{x \in \mathbb{R}^n; x \neq 0} \frac{|Ax|}{|x|} \tag{2.23}$$

where $|\cdot|$ is the standard Euclidian norm.

From submultiplicativity, we have that for any $k$

$$\left\| \frac{\partial \boldsymbol{h}_k}{\partial \boldsymbol{h}_{k-1}} \right\| \leq \|\boldsymbol{W}^{\intercal}_{hh}\| \| \operatorname{diag}\left( \sigma'_{hh}(\boldsymbol{h}_{k-1}) \right) \| \leq \gamma_{\boldsymbol{W}} \gamma_{\sigma} \tag{2.24}$$

where we have defined $\gamma_{\boldsymbol{W}} = \|\boldsymbol{W}^{\intercal}_{hh}\|$ and

$$\gamma_{\sigma} := \sup_{\boldsymbol{h} \in \mathbb{R}^n} \| \operatorname{diag}\left( \sigma'_{hh}(\boldsymbol{h}) \right) \| \tag{2.25}$$

$$= \sup_{\boldsymbol{h} \in \mathbb{R}^n} \max_i \sigma'_{hh}(\boldsymbol{h})_i \qquad \text{Operator norm of diag} \tag{2.26}$$

$$= \sup_{x \in \mathbb{R}} \sigma'_{hh}(x) \qquad \qquad \sigma_{hh} \text{ acts elementwise} \tag{2.27}$$

Substituting back into section 2.2.5, we find that

$$\left\| \frac{\partial \boldsymbol{h}_t}{\partial \boldsymbol{h}_k} \right\| = \left\| \prod_{t \geq i > k} \frac{\partial \boldsymbol{h}_i}{\partial \boldsymbol{h}_{i-1}} \right\| \leq \prod_{t \geq i > k} \left\| \frac{\partial \boldsymbol{h}_i}{\partial \boldsymbol{h}_{i-1}} \right\| \leq (\gamma_{\boldsymbol{W}} \gamma_{\sigma})^{t-k} \tag{2.28}$$

Hence, we see that a sufficient condition for vanishing gradients is for $\gamma_{\boldsymbol{W}} \gamma_{\sigma} < 1$, in which case $\left\| \frac{\partial \boldsymbol{h}_t}{\partial \boldsymbol{h}_k} \right\| \to 0$ exponentially for long timespans $t \gg k$. $\quad \square$

If $\gamma_{\sigma}$ is bounded, sufficient conditions for vanishing gradients to occur may be written as

$$\gamma_{\boldsymbol{W}} < \frac{1}{\gamma_{\sigma}} \tag{2.29}$$

1 This is true for commonly used activation functions (e.g. $\gamma_\sigma = 1$ for $\sigma_{hh} = \tanh$, $\gamma_\sigma = 0.25$ for

2 $\sigma_{hh} = $ sigmoid).

3     The converse of the proof implies that $\|W_{hh}^\top\| \geq \frac{1}{\gamma_\sigma}$ are necessary conditions for $\gamma_W \gamma_\sigma > 1$

4 and exploding gradients to occur.

## 2.2.6    Long short term memory: solving the vanishing gradient

6 In order to build a model which learns long range dependencies, vanishing gradients must be

7 avoided. This requires us to design our memory cells holding the hidden state $h$ such that

8 eq. (2.29) does not hold.

9     The *long short term memory* (LSTM) memory cell was proposed by Hochreiter and Schmid-

10 huber [14] as a method for dealing with the vanishing gradient problem. It does so by enforcing

11 *constant error flow* on section 2.2.5, that is

$$W_{hh}^\top \sigma'_{hh}(h_t) = I \tag{2.30}$$

13 where $I$ is the identity matrix.

14     Integrating the above differential equation yields $W_{hh}\sigma_{hh}(h_t) = h_t$. Since this should hold

15 for any hidden state $h_t$, this means that:

16     1. $W_{hh}$ must be full rank

17     2. $\sigma_{hh}$ must be linear

18     3. $W_{hh} \circ \sigma_{hh} = I$

19     In the *constant error carousel* (CEC), this is ensured by setting $\sigma_{hh} = W_{hh} = I$. This may

20 be interpreted as removing time dynamics on $h$ in order to permit error signals to be transferred

21 backwards in time (section 2.2.5) without modification (i.e. $\forall t \geq k : \frac{\partial h_t}{\partial h_k} = I$).

22     In addition to using a CEC, a LSTM introduces three gates controlling access to the CEC:

23     • **Input gate**: scales input $x_t$ elementwise by $i_t \in [0, 1]$, *writes* to $h_t$

24     • **Output gate**: scales output $y_t$ elementwise by $o_t \in [0, 1]$, *reads* from $h_t$

25     • **Forget gate**: scales previous cell value $h_{t-1}$ by $f_t \in [0, 1]$, *resets* $h_t$

Mathematically, the LSTM model is defined by the following set of equations:

$$i_t \;=\; \text{sigmoid}(W_{xi}x_t + W_{yi}y_{t-1} + b_i) \tag{2.31}$$

$$o_t \;=\; \text{sigmoid}(W_{xo}x_t + W_{yo}y_{t-1} + b_o) \tag{2.32}$$

$$f_t \;=\; \text{sigmoid}(W_{xf}x_t + W_{yf}y_{t-1} + b_f) \tag{2.33}$$

$$h_t \;=\; f_t \odot h_{t-1} + i_t \odot \tanh(W_{xh}x_t + y_{t-1}W_{yh} + b_h) \tag{2.34}$$

$$y_t \;=\; o_t \odot \tanh(h_t) \tag{2.35}$$

where $\odot$ denotes elementwise multiplication of vectors.

Notice that the gates ($i_t$, $o_t$, and $f_t$) controlling flow in and out of the CEC are all time varying. This can be interpreted as a mechanism enabling LSTM to explicitly learn which error signals to trap in the CEC and when to release them [14], enabling error signals to potentially be transported across long time lags.
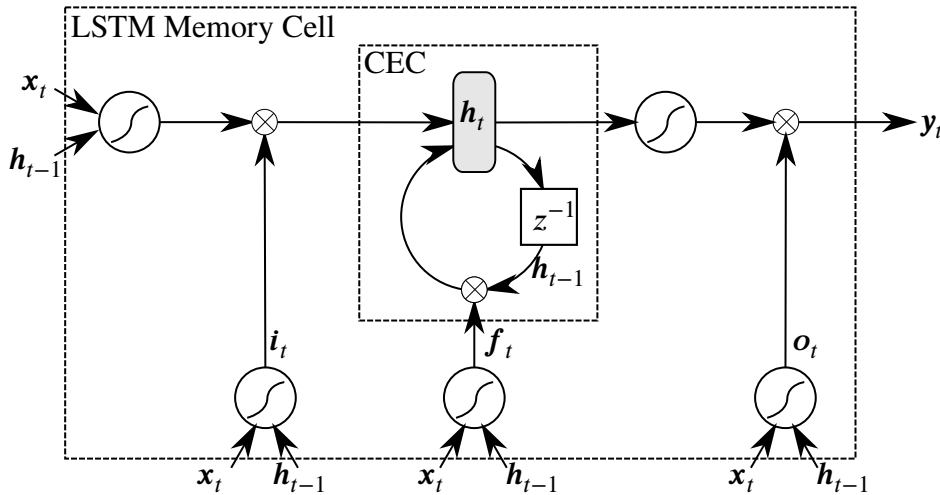


Fig. 2.13 Schematic for a single LSTM memory cell. Notice how the gates $i_t$, $o_t$, and $f_t$ control access to the constant error carousel (CEC).

Some authors define LSTMs such that $h_t$ is not used to compute gate activations, referring to the case where $h_t$ is connected as "peephole connections"[11]. We will use LSTM to refer to the system of equations as written above.

**Practicalities for successful applications of LSTMs**

Many applications of LSTMs

> fliang: cite examples

share some common practical techniques for ensuring successful training. Perhaps most important is *gradient norm clipping* [18, 20] where the gradient is rescaled whenever it exceeds

a threshold. This is necessary because while vanishing gradients are mitigated by the use of CECs, LSTMs do not explicitly protect against exploding gradients.

Another common practice is the use of methods for reducing overfitting and improving generalization. In particular, *dropout* [25] can be applied to the connections between memory cells in a stacked RNN to regularize the learned features to be more robust to noise [29]. Additionally, *batch normalization*[15] can also be applied to to the memory cell hidden states to reduce co-variate shifts, accelerate training, and improve generalization.

Finally, applications of RNNs to long sequences can incur a prohibitively high cost for a single parameter update[Sutskever]. For instance, computing the gradient of an RNN on a sequence of length 1000 costs the equivalent of a forward and backward pass on a 1000 layer feedforward network. This issue is typically addressed by only back-propogating error signals a fixed number of timesteps back in the unrolled network, a technique known as *truncated BPTT*[28]. As the hidden states in the unrolled network have nevertheless been exposed to many timesteps, learning of long range structure is still possible.

# References

[1] Bengio, Y. and Delalleau, O. (2011). On the expressive power of deep architectures. *Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics)*, 6925 LNAI:18–36.

[2] Bengio, Y., Simard, P., and Frasconi, P. (1994). Learning Long-Term Dependencies with Gradient Descent is Difficult. *IEEE Transactions on Neural Networks*, 5(2):157–166.

[3] Bharucha, J. J. and Todd, P. M. (1989). Modeling the perception of tonal structure with neural nets. *Computer Music Journal*, 13(4):44–53.

[4] Cooper, G. and Meyer, L. B. (1963). *The rhythmic structure of music*, volume 118. University of Chicago Press.

[5] Cybenko, G. (1993). Degree of approximation by superpositions of a sigmoidal function. *Approximation Theory and its Applications*, 9(3):17–28.

[6] Denton, C. and Fillion, M. (1997). The history of musical tuning and temperament during the classical and romantic periods.

[7] Eck, D. and Lapalme, J. (2008). Learning musical structure directly from sequences of music. *University of Montreal, Department of Computer Science, CP*, 6128.

[8] Elman, J. L. (1990). Finding structure in time. *Cognitive science*, 14(2):179–211.

[9] Franklin, J. A. (2006). Recurrent neural networks for music computation. *INFORMS Journal on Computing*, 18(3):321–338.

[10] Freedman, D. (2015). *Correlational Harmonic Metrics: Bridging Computational and Human Notions of Musical Harmony*. PhD thesis.

[11] Gers, F. A. and Schmidhuber, J. (2000). Recurrent nets that time and count. In *Neural Networks, 2000. IJCNN 2000, Proceedings of the IEEE-INNS-ENNS International Joint Conference on*, volume 3, pages 189–194. IEEE.

[12] Goller, C. and Kuchler, A. (1996). Learning task-dependent distributed representations by backpropagation through structure. In *Neural Networks, 1996., IEEE International Conference on*, volume 1, pages 347–352. IEEE.

[13] Handel, S. (1993). *Listening: An introduction to the perception of auditory events*. The MIT Press.

[14] Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8):1735–1780.

[15] Ioffe, S. and Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*.

[16] Jordan, M. I. (1997). Serial order: A parallel distributed processing approach. *Advances in psychology*, 121:471–495.

[17] Linnainmaa, S. (1970). The representation of the cumulative rounding error of an algorithm as a taylor expansion of the local rounding errors. *Master's Thesis (in Finnish), Univ. Helsinki*, pages 6–7.

[18] Mikolov, T. (2012). *Statistical Language Models Based on Neural Networks*. PhD thesis.

[19] Papadopoulos, G. and Wiggins, G. (1999). Ai methods for algorithmic composition: A survey, a critical view and future prospects. In *AISB Symposium on Musical Creativity*, pages 110–117. Edinburgh, UK.

[20] Pascanu, R., Mikolov, T., and Bengio, Y. (2012). On the difficulty of training recurrent neural networks. *Proceedings of The 30th International Conference on Machine Learning*, (2):1310–1318.

[21] Pearce, M., Meredith, D., and Wiggins, G. (2002). Motivations and methodologies for automation of the compositional process. *Musicae Scientiae*, 6(2):119–147.

[22] Ramage, D. (2007). Hidden markov models fundamentals. *Lecture Notes. http://cs229. stanford. edu/section/cs229-hmm. pdf*.

[23] Randel, D. M. (1999). *The Harvard concise dictionary of music and musicians*. Harvard University Press.

[24] Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1988). Learning representations by back-propagating errors. *Cognitive modeling*, 5(3):1.

[25] Srivastava, N., Hinton, G. E., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2014). Dropout: a simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(1):1929–1958.

[Sutskever] Sutskever, I. Training Recurrent Neural Networks - Ilia Sutskever - PhD thesis.

[27] Terhardt, E. (1974). Pitch, consonance, and harmony. *The Journal of the Acoustical Society of America*, 55(5):1061–1069.

[28] Williams, R. J. and Peng, J. (1990). An efficient gradient-based algorithm for on-line training of recurrent network trajectories. *Neural computation*, 2(4):490–501.

[29] Zaremba, W., Sutskever, I., and Vinyals, O. (2014). Recurrent neural network regularization. *arXiv preprint arXiv:1409.2329*.