



**VYSOKÉ UČENÍ TECHNICKÉ V BRNĚ**

BRNO UNIVERSITY OF TECHNOLOGY

**FAKULTA INFORMAČNÍCH TECHNOLOGIÍ**

FACULTY OF INFORMATION TECHNOLOGY

**ÚSTAV POČÍTAČOVÉ GRAFIKY A MULTIMÉDIÍ**

DEPARTMENT OF COMPUTER GRAPHICS AND MULTIMEDIA

**SEPARACE MLUVČÍCH V ČASOVÉ DOMÉNĚ**

TIME DOMAIN AUDIO SEPARATION

**BAKALÁŘSKÁ PRÁCE**

BACHELOR'S THESIS

**AUTOR PRÁCE**

AUTHOR

**JIŘÍ PEŠKA**

**VEDOUcí PRÁCE**

SUPERVISOR

**ing. KATEŘINA ŽMOLÍKOVÁ,**

**BRNO 2020**

## Zadání bakalářské práce



23185

Student: **Peška Jiří**

Program: Informační technologie

Název: **Separace mluvčích v časové doméně pomocí neuronové sítě**  
**Time-Domain Neural Network Based Speaker Separation**

Kategorie: Zpracování řeči a přirozeného jazyka

Zadání:

1. Seznamte se s problémem separace mluvčích pomocí neuronových sítí.
2. Seznamte se s metodou TasNet pro jednokanálovou separaci signálu v časové doméně.
3. Implementujte danou metodu s využitím vhodného toolkitu (např. PyTorch, Keras).
4. Otestujte systém na vhodném datasetu. Zaměřte se na vyhodnocení vlivu velikosti sítě na přesnost.
5. Navrhněte a diskutujte možné zlepšení použité metody.

Literatura:

- Luo, Yi, and Nima Mesgarani. "Conv-TasNet: Surpassing Ideal Time-Frequency Magnitude Masking for Speech Separation." *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 27.8 (2019): 1256-1266.
- dle doporučení vedoucího

Podrobné závazné pokyny pro vypracování práce viz <https://www.fit.vut.cz/study/theses/>

Vedoucí práce: **Žmolíková Kateřina, Ing.**

Vedoucí ústavu: Černocký Jan, doc. Dr. Ing.

Datum zadání: 1. listopadu 2019

Datum odevzdání: 14. května 2020

Datum schválení: 5. listopadu 2019

## Abstrakt

Práce se zabývá využitím konvolučních neuronových sítí pro automatickou separaci mluvčích v akustickém prostředí. Cílem je implementovat neuronovou síť podle architektury TasNet za použití frameworku PyTorch, natrénovat síť s různými hodnotami hyperparametrů a porovnat kvalitu separací vzhledem k velikosti sítě.

Architektura oproti dosavadním metodám, které převáděly vstupní směs do časově-frekvenční reprezentace, používá konvoluční autoenkodér, který vstupní směs převádí do nezáporné reprezentace, která je optimalizovaná pro extrakci jednotlivých mluvčích. Samotné separace je docíleno aplikací masek, které jsou odhadnuty v separačním modulu. Modul tvoří opakující se posloupnost konvolučních bloků se zvyšující se dilatací, která napomáhá k modelování časových závislostí ve zpracovávané směsi.

K vyhodnocení přesnosti bylo použita metrika Signal to Distortion Ratio (SDR), která určuje poměr zastoupení šumu a zvuku v nahrávce. Vyhodnocování proběhlo na množině směsí dvou mluvčích, která byla vygenerována náhodnou kombinací mluvčích z Wall Street Journal datasetu (WSJ0). Natrénováním několika modelů s různými hodnotami hyperparametrů bylo možno pozorovat závislost mezi velikostí sítě a hodnotou SDR. Zatímco menší síť dosahovala, po X epochách trénování, přesnosti XY, větší síť dosahovala až XX.

**[[Doplnit SDR přesnost do odstavce vyše. Jeste 4. cast?]]**

## Abstract

**[[Preložit do anglickiny CZ abstrakt]]** Thesis is about usage of convolutional neural networks for automational speech separation in acoustic environment. The goal is to implement neural network by following TasNet architecture in pytorch framework, train this network with various values of hyper parameters and to compare the quality of separation based on the size of the network.

In contrast to older architectures, that transformed input mixture into time–frequency representation, this architecture use convolutional autoencoder, which transforms input mixture into non–negative representation optimized for speaker extraction. Separation is achieved by applying the masks, which are estimated in the separation module. This module consists of stacked convolutional blocks with increasing dilation, which helps with modelling of long–term time dependencies in the processed speech.

For evaluation of preci. Ye.

## Klíčová slova

neuronové sítě, zpracování řeči, konvoluční neuronová síť, autoenkodér, separace mluvčích, strojové učení, tasnet, feed forward, hluboké učení

## Keywords

artificial neural networks, speech processing, convolutional neural networks, autoencoder, speech separation, machine learning, tasnet, feed forward, deep learning

## Citace

PEŠKA, Jiří. *Separace mluvčích v časové doméně*. Brno, 2020. Bakalářská práce. Vysoké učení technické v Brně, Fakulta informačních technologií. Vedoucí práce ing. Kateřina Žmolíková,

# Separace mluvčích v časové doméně

## Prohlášení

Prohlašuji, že jsem tuto bakalářskou práci vypracoval samostatně pod vedením ing. Kateřiny Žmolíkové. Další informace mi poskytli... Uvedl jsem všechny literární prameny, publikace a další zdroje, ze kterých jsem čerpal.

.....

Jiří Peška

10. července 2020

## Poděkování

V této sekci je možno uvést poděkování vedoucímu práce a těm, kteří poskytli odbornou pomoc (externí zadavatel, konzultant apod.).

# Obsah

<b>1</b>	<b>Úvod</b>	<b>3</b>
<b>2</b>	<b>Separace mluvcích</b>	<b>4</b>
2.1	Koktejl párty efekt . . . . .	4
2.2	Metody pro separaci jednokanálových nahrávek . . . . .	4
2.2.1	CASA – computational auditory scene analysis . . . . .	4
2.2.2	NMF – non-negative matrix factorization . . . . .	5
2.2.3	Hluboké neuronové sítě . . . . .	5
2.3	Vyhodnocovací metriky . . . . .	6
2.3.1	Source to distortion ratio . . . . .	6
2.3.2	Perceptual evaluation of speech quality . . . . .	6
2.3.3	Short-time objective intelligibility . . . . .	7
<b>3</b>	<b>Neuronové sítě</b>	<b>8</b>
3.1	Organizace feedforward sítí . . . . .	8
3.2	Umělý neuron . . . . .	10
3.3	Aktivační funkce . . . . .	11
3.4	Konvoluční vrstvy . . . . .	14
3.4.1	Bottle-neck vrstva . . . . .	17
3.5	Učení neuronových sítí . . . . .	17
3.5.1	Objektivní funkce . . . . .	17
3.5.2	Backpropagation . . . . .	19
3.5.3	Underfitting a generalizace . . . . .	20
3.5.4	Reziduální spojení . . . . .	21
3.5.5	Význam validační množiny v trénování . . . . .	22
<b>4</b>	<b>TasNet - Time-Domain Audio Separation Network</b>	<b>23</b>
4.1	Konvoluční auto-enzodér . . . . .	24
4.1.1	Enkódování směsi . . . . .	24
4.1.2	Dekódování extrahovaných mluvcích . . . . .	24
4.2	Separací modul . . . . .	24
4.2.1	Temporal Convolutional Network . . . . .	25
4.2.2	Konvoluční bloky - možná ne jako podkapitolu . . . . .	27
4.3	Normalizace . . . . .	28
<b>5</b>	<b>Implementace sítě</b>	<b>30</b>
5.1	Implementace tříd a modelu . . . . .	31
5.2	Segmentace nahrávek . . . . .	33

5.3	Pomocné skripty . . . . .	36
<b>6</b>	<b>Experimenty a vyhodnocení</b>	<b>38</b>
6.1	Dataset . . . . .	38
6.2	Průběh trénování . . . . .	39
6.3	Experimenty s modely . . . . .	40
6.3.1	Výsledky referenčního modelu TasNet . . . . .	41
6.3.2	Vliv hyper-parametru X . . . . .	41
6.3.3	Vliv hyper-parametru R . . . . .	41
6.3.4	Vliv hyper-parametru délka-segmentu . . . . .	41
6.3.5	Vyhodnocení a porovnání modelů . . . . .	41
6.4	Možná rozšíření a navrhnutá vylepšení . . . . .	41
<b>7</b>	<b>Závěr</b>	<b>42</b>
	<b>Literatura</b>	<b>43</b>

# Kapitola 1

## Úvod

Zpracování řeči hraje v dnešní době důležitou roli v mnoha rozličných oborech. Mezi jeden z hlavních úkolů bezesporu patří separace zdrojů v zaznamenaném signálu, který může být složen ze signálů  $N$  mluvčích, ale i nechtěného hluku okolí. Vyřešení problému je předpoklad k dalším úkonům jako identifikace konkrétního mluvčího nebo třeba přepis konverzace na text. Se stále se zrychlujícím vývojem počítačů a s jejich zvyšujícím se výkonem se do popředí dostávají metody zpracování řeči založené na neuronových sítích, které v mnoha ohledech předčí ostatní algoritmy strojového učení.

Separace mluvčích v časové doméně dosahuje mimořádných výsledků v porovnání s dosavadními metodami založenými na převodu signálu z časové domény do frekvenční domény pomocí algoritmu STFT (Short-Time Fourier Transform). Taková reprezentace signálu není optimalizovaná pro separaci řeči a nemusí pro tento úkol být optimální. V architektuře, která je navržena v referenční studii s názvem **TasNet: Surpassing Ideal Time-Frequency Masking for Speech Separation** [12], je vstupní signál převeden do nezáporné reprezentace, která je optimální pro extrakci jednotlivých mluvčích. Silnou stránkou systému je hluboká architektura sítě, která lépe modeluje dlouhodobé závislosti v signálu. Zároveň se ale musí vypořádat s problémy, které hluboké neuronové sítě mohou přinášet.

Úkolem této práce je nastudovat si problematiku neuronových sítí a jejich základní principy, seznámit se problémem separace mluvčích pomocí neuronových sítí a následně implementovat síť podle architektury TasNet pro separaci mluvčích v časové doméně, která byla navržena a popsána ve studii [12]. Poté tuto neuronovou síť natrénovat s různými kombinacemi hodnot hyperparametrů, které ovlivňují velikost sítě a její vlastnosti, a nakonec porovnat přesnost a kvalitu separace mezi jednotlivými, různě velkými sítěmi a s výsledky, kterých bylo dosaženo ve studii. Přesnost a kvalita separace bude vypočítána pomocí metrik určených k hodnocení kvality řečového signálu. Sítě budou testovány a vyhodnocovány na testovací množině jednokanálových směsí dvou mluvčích.

Text je rozdělen na několik kapitol, které postupně pokrývají různé logické celky. V kapitole 2 je popsán problém separace mluvčích a dosavadní přístupy k jeho řešení. Kapitola 3 se zabývá základními vlastnostmi a principy neuronových sítí. Kapitola 4 je věnována architektuře TasNet, kterou se moje práce zabývá. Dále kapitola 5 popisuje implementaci sítě a pomocných nástrojů a nakonec v kapitole 6 jsou experimenty a vyhodnocování přesnosti sítě.

## Kapitola 2

# Separace mluvčích

Živočichové mají vrozenou schopnost zaměřit se na jeden konkrétní zvuk, zatímco všechny ostatní dokážou potlačit. Naučit ale této dovednosti počítač se ukázalo jako obtížný úkol, který je překážkou ASR (*automatic speech recognition*) systémům v rozsáhlejší nasazení do běžného života. Tento problém se nazývá koktejl párty. Vyřešení tohoto problému, tedy separování zdrojů, umožňuje širokou aplikaci v oblasti zpracování řeči, jako přepis konverzací na text, ovládání počítače hlasem, používání hlasových asistentů či identifikace konkrétního mluvčího v prostředí s mnoha současně mluvícími lidmi.

### 2.1 Koktejl párty efekt

Tento problém se nejlépe dá popsat jako schopnost soustředit se na jeden konkrétní zdroj zvuku z mnoha souběžných zdrojů jako hudba, konkrétní hudební nástroj, řeč dalších lidí či okolní zvuky, a přepínat mezi nimi dynamicky pozornost.

**[[vzniklo pro jeho vyreseni hafo altofirmu, ktere dosahovaly nejakyh vysledku. jeho vyresenim by se mohlo pouzivat asr na prepis reci an text a pod. Nejlepsich vysledku dosahly hluboke neuronove site.]]**

### 2.2 Metody pro separaci jednokanálových nahrávek

Pro separaci zdrojových signálů z nahrávek směsí existuje mnoho algoritmů založených na zpracování signálů, z nichž asi nejznámější jsou computational auditory scene analysis (CASA) [\[\[ref?\]\]](#) a non-negative matrix factorization (NMF) [\[\[ref?\]\]](#), používané pro separaci jednokanálových nahrávek. V posledních letech se do popředí dostaly přístupy založené na hlubokých neuronových sítích, které předčily svou přesností a výkonností dosavadní algoritmy. Pro separaci vícekanálových nahrávek se používají metody *beamforming* nebo state-of-the-art *Multi-channel blind source separation*, které ale jsou nad rámec této práce. Tato kapitola shrnuje některé metody pro řešení separace, které jsou zmíněny v článku [\[16\]](#).

#### 2.2.1 CASA – computational auditory scene analysis

Metoda CASA je založena na procesu v lidském mozku, který separuje zdroje ze směsi mluvčích a simuluje vysokoúrovňové chování lidského sluchu. Pro separaci jsou většinou manuálně navrhnutá segmentační pravidla pro operování nad nízkourovňovými příznaky k odhadu časově frekvenční ( $T$ - $F$ ) masky, která izoluje komponenty signálu jednotlivých



mluvčích a následně je použita pro rekonstrukci všech zdrojů. **[[Hodit sem teda ty zdroje jako v 3.1 pro ASA, CASA a pod? jak jsou tam asi 3 po sobe?]]**

Přestože tato metoda byla i nadále rozvíjena, tak má mnoho nevýhod jako špatnou generalizaci v důsledku manuálního vytváření pravidel, nemožnost automatického učení z dat nebo nemožnost použití na separaci jiných zdrojů než je lidská řeč, které značně omezují její použití v mnoha reálných případech.

### 2.2.2 NMF – non-negative matrix factorization

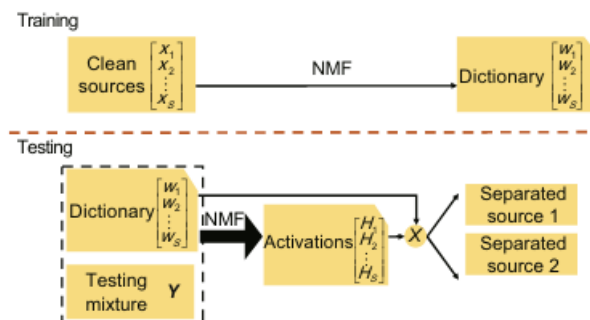
Nezáporná maticová faktorizace (NMF) **[[[(Lee and Seung, 2001)]]]** patří mezi daty řízené metody a je založena na předpokladu, že struktura spektrogramu může být reprezentována malým počtem bází.

Pro NMF platí

$$Y = \sum s W_s H_s \quad (2.1)$$

kde každý zdroj  $s$  je modelován aproximací nízkého řádu nezápornými maticemi  $W_s$ , která reprezentuje slovník a  $H_s$ , která reprezentuje aktivační funkci. Ty jsou sečteny a dohromady tvoří výslednou směs  $Y$ .

Metoda prochází trénovací a testovací fází. V trénovací fázi je každý zdroj dekomponován a mapován na množinu bází a aktivací a tím je zformován slovník  $W$  pro tento zdroj. V testovací fázi jsou naučené slovníky jednotlivých zdrojů  $s$  spojeny do jednoho, který dále není upravován. Následuje optimalizace aktivační matice  $H$  pro každý zdroj. Každý zdroj je následně ze směsi rekonstruován s použitím výsledných bází a aktivací.



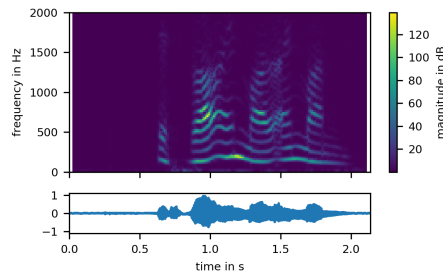
Obrázek 2.1: Znázornění trénovací a testovací fáze. Obrázek byl inspirován [16]

Stejně jako metoda CASA, i tato má mnoho nevýhod, které omezují její použití pro řešení reálných problémů separace.

### 2.2.3 Hluboké neuronové sítě

Nejlépe si při řešení problému koktejl párty vedou postupy založené na hlubokém učení a neuronových sítích. Metody převáděly směs do časově-frekvenční reprezentace znázorněné na obrázku 2.2 pomocí STFT (*short time fourier transformation*).

**[[rozepsat]]**



Obrázek 2.2: Časově–frekvenční reprezentace řeči, neboli spektrum

## 2.3 Vyhodnocovací metriky

Pro vyhodnocení kvality separace se používají různé metriky. Jejich volba závisí na řešeném problému. Jejich výpočet probíhá na testovacím datasetu. Nejčastěji používané metriky pro měření kvality separace jsou metriky pro vyhodnocování slepé separace zdrojů (BSS), dále short-time objective intelligibility (STOI) a procentuální vyhodnocení kvality řeči (PESQ). Metriky pro BSS jsou source-to-distortion ratio (SDR), signal-to-inference ratio (SIR) a signal-to-artifacts ratio (SAR).

### 2.3.1 Source to distortion ratio

SDR [22] je jedna z nejznámějších metrik pro měření kvality separace mluvčích. Hodnota je udávána v decibelech (dB). Na výstupu separačního algoritmu se může nacházet rušení, které lze kategorizovat jako inference a artefakty.

Inference odkazuje na přetrvávající podíl druhého zdroje na separované nahrávce prvního zdroje. SIR je tedy metrika zachycující schopnost algoritmu odstranit ostatní zdroje a zachovat pouze ten požadovaný.

Během separace může algoritmus do separovaných nahrávek zanechat artefakty, které se mohou projevit jako rušení či zvuky, které na původní nahrávce nebyly. SAR je metrika zachycující schopnost algoritmu separovat nahrávky bez toho, aby do nich zanesl nechtěné artefakty [21].

Metriky SIR a SAR lze zkombinovat do jediné metriky SDR, která udává celkovou kvalitu separace.

Hodnota SDR se vypočítá jako

$$SDR = 10 * \log_{10} \frac{\|s_{target}\|^2}{\|e_{interf} + e_{noise} + e_{artif}\|^2} \quad (2.2)$$

[[popsat einfer atd nevím jestli ANO nebo NE]]

### 2.3.2 Perceptual evaluation of speech quality

Procentuální vyhodnocení kvality řeči (PESQ) [17] je standardizovaná objektivní metoda pro vyhodnocování kvality řeči na end-to-end síti. PESQ byla vyvinuta pro předpověď subjektivního skóre *mean opinion score* (MOS), kde posluchač hodnotil kvalitu připojení (řeči) jednou z pěti možností (bad, poor, fair, good, excellent). Těmto možnostem odpovídaly číselné hodnoty [[1–5 udelat správně pomlčku]], které pak byly zprůměrovány, což nakonec byla finální hodnota MOS.

Hodnota PESQ se pohybuje od  $-0.5$  do  $4.5$  s tím, že vyšší hodnota znamená lepší procentuální kvalitu [?]. Hodnota je určena na základě porovnání originálního signálu a přeneseného signálu skrze komunikační systém. Lze ho ale aplikovat i pro vyhodnocení separace mluvčích, kde referenční nahrávka je porovnána s nahrávkou, která byla vyprodukována neuronovou sítí.

**[[popsat víc nebo ne.. treba to ze je tam nejake alignment, input filter, FFT viz odkaz]]**

### 2.3.3 Short-time objective intelligibility

Short-time objective intelligibility (STOI) [19] je metrika založená na hodnotě korelačního koeficientu mezi časově-frekvenčním spektrem referenčního a odhadnutého signálu řeči v krátkých překrývajících se segmentech [20]. Spektrum čisté a separované nahrávky je získáno pomocí short-time fourier transform (STFT).

Oba signály jsou dekomponovány pomocí  $\frac{1}{3}$ -oktávových filtrů pomocí STFT, následně segmentovány na krátké úseky s 50% překrytím, normalizovány, klipnuty a nakonec porovnány průměry hodnot korelačních koeficientů [19]. Hodnota STOI se pohybuje v rozmezí 0 a 1, korespondující s procentuální úspěšností.

Bylo navrženo mnoho konvenčních přístupů a metod pro řešení problému koktejl párty a s ním spojené separace mluvčích, ale většina z nich nedosahuje dostatečné přesnosti nebo výkonu k použití v reálných situacích. Přístupy založené na hlubokém učení a neuronových sítích předčily dosavadní techniky jak v přesnosti, tak ve výkonu, což umožňuje jejich nasazení v běžném životě.

## Kapitola 3

# Neuronové sítě

V dnešní době zažívají neuronové sítě díky výkonosti počítačů velký rozmach. Jejich využití prostupuje skrze mnohé vědní obory a dokáží řešit celou řadu problémů, ve kterých dosahují výborných výsledků, které zdaleka předčily dosavadní postupy.

Neuronové sítě (*artificial neural networks*) jsou výpočetní model, který je inspirovaný strukturou lidského mozku, ve kterém je obrovské množství propojených a komunikujících neuronů. Ty se skládají ze vstupních dendridů, výstupních axonů a samotného těla neuronu. Na základě vnitřního potenciálu a vstupních hodnot je po přesažení prahové hodnoty neuron vybuděn a je vyslán signál na výstupní axon. Signál je nakonec předán dalším neuronům skrze jejich vstupní dendridy [9, p. 65–66].

Účelem neuronové sítě je naučit se plnit zadanou úlohu. Rozdíl oproti běžným algoritmům je ale ten, že způsob, jakým síť má problém řešit, není explicitně naprogramován, ale je postupně naučen. Základní způsoby učení jsou s učitelem (*supervised*) a bez učitele (*unsupervised*).

**Učení s učitelem**, pod které spadá i tato práce, spočívá v mapování vstupních dat na data výstupní na základě vzorových příkladů dvojic vstup–výstup. Množině takových dvojic se říká trénovací dataset. Další vlastností trénovacích dat určených pro učení s učitelem, je jejich označení (*label*).

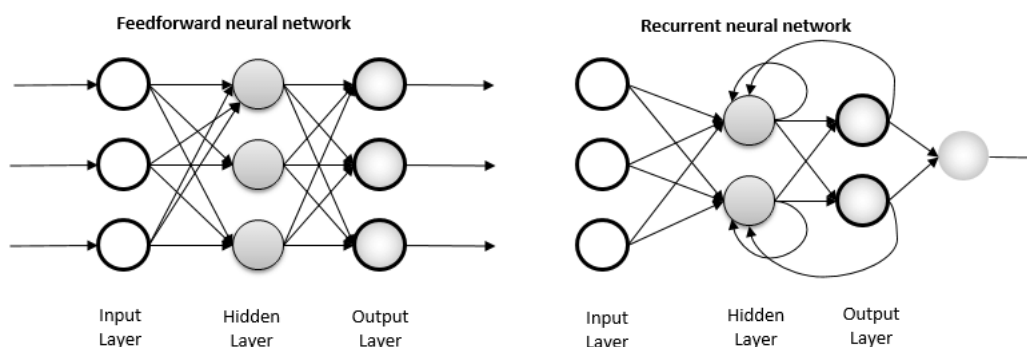
Mezi problémy, které se dají řešit neuronovými sítěmi patří klasifikační a regresní problémy. Konkrétní příklad z oblasti klasifikace může být rozpoznávání objektů na obraze, psaného písma nebo detekce obličejů na videu, ale i mnohé aplikace ve zpracování řeči. Některé z problémů lze řešit sítí, která se skládá z jednotek neuronů, zatímco složitější problémy vyžadují mnohem větší kapacitu sítě a ta pak může obsahovat i tisíce neuronů.

Upravovat samotnou strukturu a chování neuronové sítě lze pomocí jejích hyperparametrů, což jsou parametry určující nastavení neuronové sítě a trénovacího algoritmu. Tyto parametry musí být určeny před začátkem trénování a většinu nelze později měnit. Mohou určovat kapacitu modelu, velikost záběrného pole, velikost filtrů, ale i regulovat samotný proces učení nastavením počtu epoch, po které se model má učit, nebo také počet dat předaný v jedné dávce (tzv. *minibatches*) během učení sítě.

### 3.1 Organizace feedforward sítí

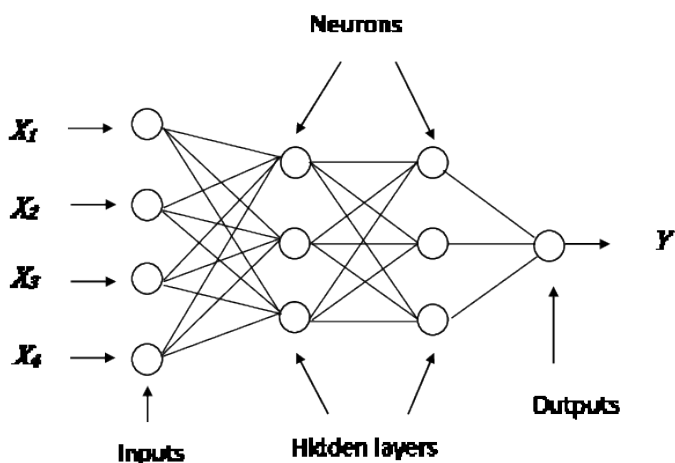
Feed forward neuronové sítě (*Multi Layer Perceptron – MLP*) jsou typ umělých neuronových sítí, kde se nevyskytují cykly ve výpočetním grafu, tedy zpětná propojení vrstev, takže informace se pohybuje pouze jedním směrem, od vstupní vrstvy přes skryté vrstvy

až po vrstvu výstupní. Sítě, které obsahují cykly, se nazývají rekurentní. Rozdíl znázorňuje obrázek 3.1.



Obrázek 3.1: Příklad grafu feed forward sítě a rekurentní neuronové sítě. Lze si všimnout orientace šipek u feed forward sítě, které směřují pouze jedním směrem, zatímco u rekurentní sítě šipky směřují i k předešlým uzlům grafu [15].

Struktura neuronové sítě je organizována do vrstev, které se skládají z neuronů. Feedforward síť je tvořena třemi typy vrstev (viz obrázky 3.1 a 3.2). Vstupní vrstva slouží k předání hodnot do sítě, ale nijak tyto hodnoty nemodifikuje. Nezměněné jsou zkopírovány do první skryté vrstvy. Následují skryté vrstvy, z nichž poslední je napojena na výstupní vrstvu. Ta má obvykle méně neuronů než předešlé vrstvy a hodnoty na výstupu mohou představovat třídy, do kterých má být klasifikován vstup v případě klasifikátorů, nebo predikce hodnot na základě vstupních dat v případě regrese. S počtem jednotlivých vrstev souvisí pojem hloubka sítě, která je rovna počtu všech vrstev neuronové sítě od vstupní až po výstupní vrstvu. Pojem "[**uvozovky lepe**]" "hluboká neuronová síť" označuje takovou síť, která má dvě nebo více skrytých vrstev. Existuje mnoho typů vrstev, například plně propojené, pooling, s přeskočením nebo vrstvy konvoluční.

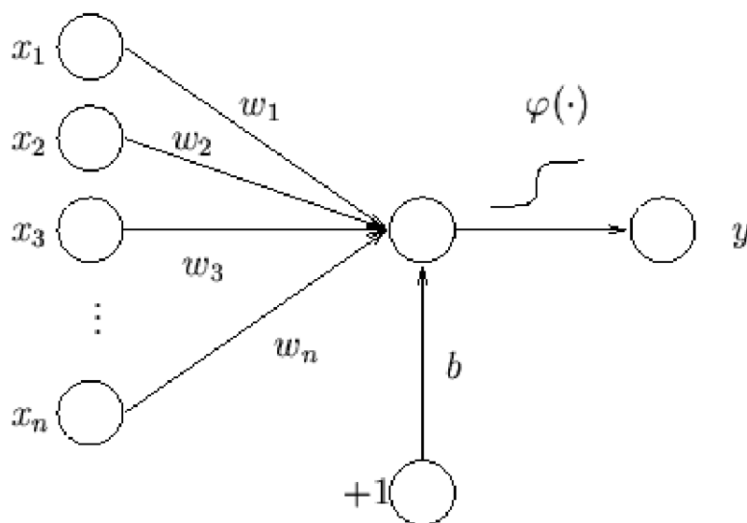


Obrázek 3.2: Schéma neuronové sítě, která má 2 skryté vrstvy

## 3.2 Umělý neuron

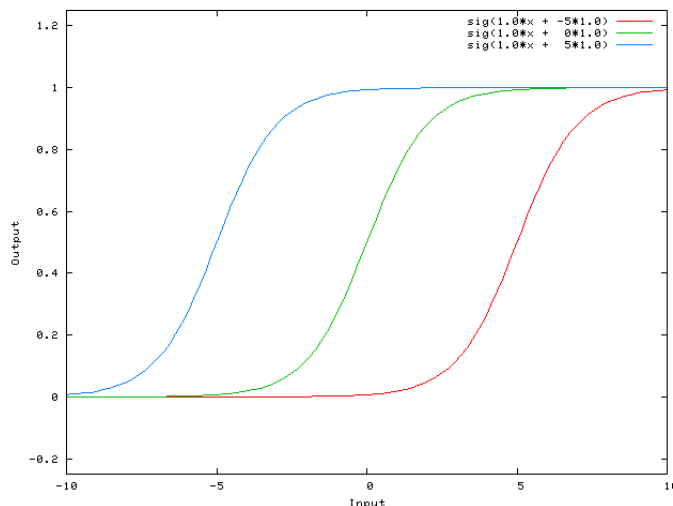
Základní stavební jednotka neuronových sítí je umělý neuron (*artificial neuron*) (viz obrázek 3.3). Tento model je založen na principu reálných neuronů, které se nacházejí v organizmech. Umělý neuron obsahuje libovolně mnoho vstupních propojení, přes které se mu předávají data v podobě vstupního vektoru  $\vec{x} = [x_1, x_2, \dots, x_n], x_n \in \mathbb{R}$ . Sám neuron obsahuje hodnotu bias  $b \in \mathbb{R}$  a vektor vah  $\vec{w} = [w_1, w_2, \dots, w_n], w_n \in \mathbb{R}$ , jenž je upravován během trénování neuronu.

Výstupní hodnota závisí na vstupních datech, aktuálním vnitřním stavu (hodnoty vah a biasu) a na zvolené aktivační funkci. Vstupní hodnoty jsou váhovány, což znamená, že každá vstupní hodnota je vynásobena s vahou na daném vstupním spojení. S použitím definovaných vektorů lze napsat, že vstupní vektor je vynásoben s vektorem vah.



Obrázek 3.3: Schéma umělého neuronu

Hodnota bias  $b$ , která je přičtena k sumě násobků vah a vstupních hodnot, je prahová hodnota modifikující dobu, kdy se aktivuje neuron a změní svůj výstup. Matematicky to znamená, že s grafem aktivační funkce horizontálně pohybuje doleva nebo doprava v závislosti na tom, je-li hodnota biasu pozitivní nebo negativní. Toto posunutí je znázorněno na obrázku 3.4. V závislosti na řešeném problému může být žádoucí, aby i hodnota bias byla modifikována během učení společně s ostatními váhami.



Obrázek 3.4: Vliv hodnoty bias na aktivační funkci

Výstup neuronu se vypočítá jako

$$y = f\left(\left(\sum_{k=1}^n w_k x_k\right) + b\right) \quad (3.1)$$

kde  $f$  je aktivační funkce,  $x_k \in \mathbb{R}$  je vstupní hodnota,  $w_k \in \mathbb{R}$  je váha, kterou se vstupní hodnota vynásobí a  $b \in \mathbb{R}$  je hodnota bias.

### 3.3 Aktivační funkce

Aktivační, neboli prahová funkce, určuje výstupní hodnotu neuronu. Funkce se vybírá na základě problému, který se má neuronová síť naučit řešit. Správná volba aktivační funkce vede k lepší konvergenci učení sítě. Naopak špatná volba může vést ke stále větší odchylce od správného řešení – může divergovat. Povaha problému může vyžadovat specifické vlastnosti aktivační funkce – lineární nebo nelineární. Pro nestandardní problémy je obvykle potřeba experimentálně zjistit, která funkce bude nejlépe vyhovovat danému problému.

Pokud by veškeré aktivační funkce v modelu byly lineární, tak celkové mapování sítě by bylo omezeno pouze na lineární mapování vstupu na výstup. Reálné problémy ale lineární obvykle nejsou a v případě pokusu modelovat takovým modelem nelineární vztahy by vedlo k velice nepřesným výsledkům, který by byl zapříčiněn podučením (*underfitting*), což znamená, že model, který se učí zakódovat nějaký vzor v datasetu, je příliš jednoduchý. Proto je potřeba zavést do modelu i nelineární aktivační funkce, které tento problém řeší [9, p. 77–78].

Z pohledu učení je také důležité, aby aktivační funkce byla diferencovatelná. To umožňuje použití trénovacích metod založených na výpočtu gradientu, jako je algoritmus backpropagation.

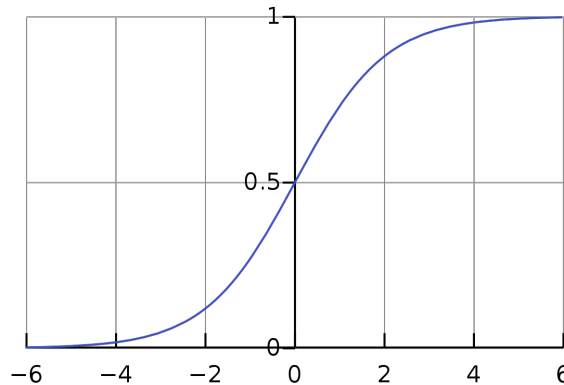
#### Logická sigmoida

Logická sigmoida je nelineární aktivační funkce používaná v hlubokém učení, nejčastěji v dopředných neuronových sítích. Graf funkce má hladký průběh, jak lze vidět na obrázku 3.5.

Je to omezená diferencovatelná funkce, definovaná jako

$$f(x) = \frac{1}{1 + \exp(-x)} \quad (3.2)$$

Používá se na výstupních vrstvách pro predikci výstupu na základě pravděpodobnosti, čehož se využívá při řešení logistických regresních problémů, u binárních klasifikátorů a dalších oblastí neuronových sítí. Podle [13] je nejčastěji používána v sítích s menším počtem vrstev (tzv. shallow networks). Mezi hlavní nevýhody patří ostrý gradient během zpětné propagace, saturace gradientu a pomalá konvergence [14].

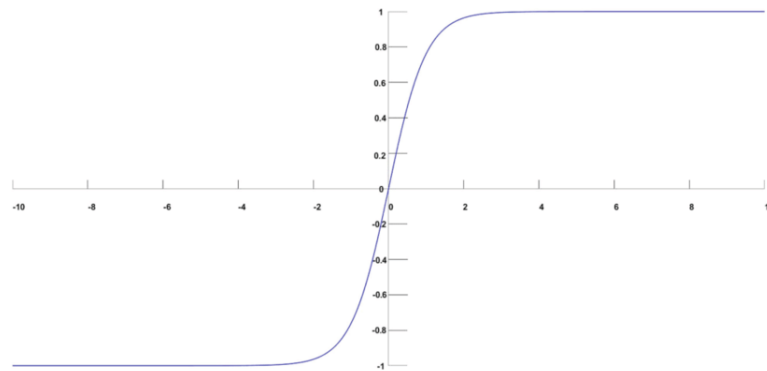


Obrázek 3.5: Graf aktivační funkce sigmoid

## Softmax

Aktivační funkce softmax [7] se často umísťuje na výstupní vrstvy a používá se pro výpočet rozdělení pravděpodobnosti z vektoru reálných čísel. Výstupem je hodnota v rozmezí 0 a 1, jak zobrazuje graf 3.6. V klasifikačních modelech výstup reprezentuje pravděpodobnosti jednotlivých tříd s tím, že cílová třída má nejvyšší hodnotu pravděpodobnosti [14]. Vypočítá se jako

$$f(x_i) = \frac{\exp(x_i)}{\sum_j \exp(x_j)} \quad (3.3)$$



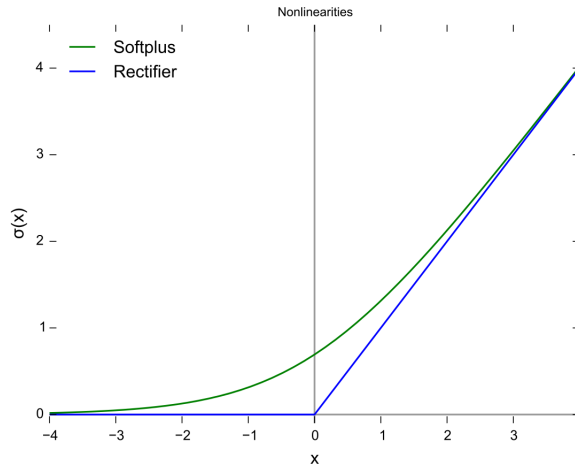
Obrázek 3.6: Graf aktivační funkce softmax[[lepsi graf udelat]]



## Rectified Linear Unit

Rectified Linear Unit (ReLU) je nejčastěji používaná aktivační funkce a stojí za nedávnými úspěchy v oblasti neuronových sítí [10]. Vyžaduje-li síť nějakou nelinearitu, je ReLU pro většinu případů ideální. Její použití zrychluje konvergenci trénovacího procesu a vede k lepším výsledkům, než u jednotek používající konvenční sigmoidní aktivační funkce [25]. Pro každou zápornou hodnotu  $x$  vrací 0 a pro kladnou hodnotu  $x$  vrací tutéž hodnotu  $x$ , jak udává rovnice

$$f(x) = \max(0, x) \quad (3.4)$$



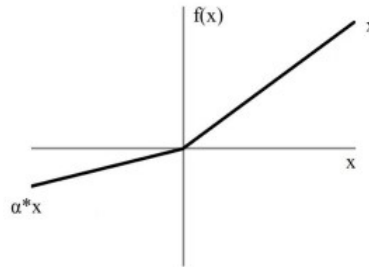
Obrázek 3.7: Graf aktivační funkce ReLU

## Parametric Rectified Linear Unit

Parametrizovaná ReLU (PReLU) [6] je nelineární aktivační funkce a je to generalizace funkce ReLU. Používá se v případě, že chceme produkovat na výstup malý nenulový gradient i v případě záporné vstupní hodnoty  $x$ , kdy je vynásobena parametrem  $\alpha$ . Parametr  $\alpha$  reguluje klesání záporné části grafu a adaptivně se učí společně s ostatními váhami během trénovacího procesu.

Předpis funkce PReLU je definován jako

$$f(x) = \begin{cases} x & \text{if } x \geq 0 \\ \alpha x & \text{if } x < 0 \end{cases} \quad (3.5)$$



Obrázek 3.8: Graf aktivační funkce PReLU

V případě, že je hodnota parametru  $\alpha = 0$ , pak se jedná o ReLU. Je-li parametr  $\alpha = c$ , kde  $c \in \mathbb{R}$  je malá konstanta (například  $c = 0.03$ ), jedná se o Leaky ReLU (LReLU) a motivací pro její zavedení bylo vyhnout se nulovému gradientu.

### 3.4 Konvoluční vrstvy

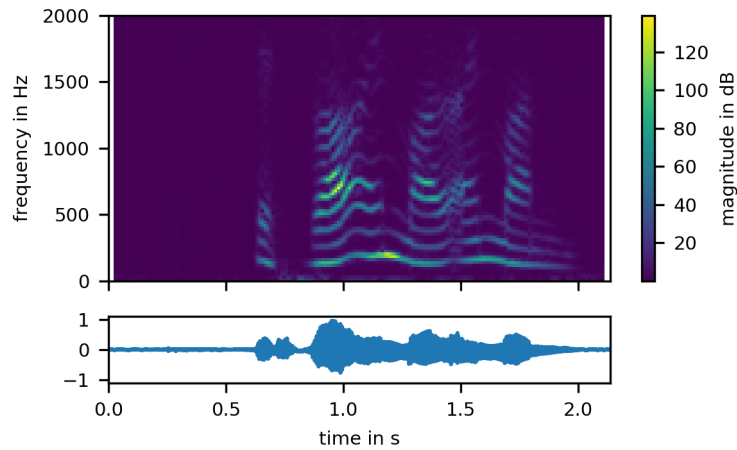
Konvoluční neuronová síť (CNN) [11] je typ feed forward neuronových sítí, která obsahuje jednu či více konvolučních vrstev. CNN jsou vhodné pro zpracování dat v mřížkovitém uspořádání. To může být ve 2-D obrázek ve formě pole pixelů či spektrum nahrávky, nebo vzorky zvukové nahrávky v čase v 1-D. Konvoluční vrstvy používají místo běžného násobení matic speciální typ lineární operace – konvoluci.

Diskrétní konvoluce je definována jako

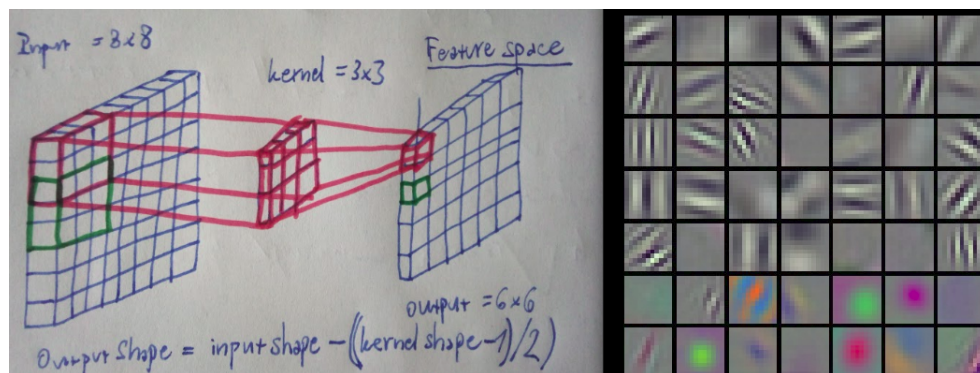
$$(f \star g)_k = \sum_{i=-\infty}^{\infty} f_i g_{k-i} = \sum_{i=-\infty}^{\infty} f_{k-i} g_i \quad (3.6)$$

kde  $\star$  je konvoluční operátor,  $f$  je funkce signálu, funkce  $g$  je konvoluční jádro a  $f_i$  a  $g_i$  jsou hodnoty funkce na indexu  $i$ .

Účelem CNN je v prvních vrstvách sítě extrahovat lokální příznaky ze vstupu. Pro 2-D vstup to mohou být úsečky v různých úhlech nebo části křivek a pro 1-D tóny o různé frekvenci. V pozdějších vrstvách neurony kombinují extrahované lokální příznaky do komplexnějších příznaků (oko, nos), až dokud nesformují v konečných vrstvách například celý obličej, znak, kočka či jiný objekt, který se síť učí detekovat. V některých modelech zpracování řeči je zvuk transformován na 2-D spektrogram (viz obrázek 3.9) a lze s ním pracovat stejně, jako při zpracování obrazu.



Obrázek 3.9: Ukázka nahrávky řeči a korespondující spektrogram, na němž jsou zobrazeny frekvence (osa y) z nahrávky řeči v čase (osa x). Zvuk se na spektrogram převádí pomocí operace short time fourier transform (STFT)



Obrázek 3.10: Obrázek vlevo zobrazuje, jak probíhá konvoluování jádra se vstupními 2-D daty a mapování lokálních příznaků na prostor příznaků (feature space), což je výstupem konvoluce. Jádro je postupně aplikováno skrze celý vstup. Obrázek napravo zobrazuje, jak mohou vypadat lokální příznaky.

Konvoluování jádra přes nějaký vstup (například obrázek či zvuk) je ekvivalentní s detektorem, který detekuje nějaký lokální příznak. Při aplikaci takového detektoru přes nějaký vstup zaznamenává všechny pozice, kde se příznak nacházel [9].

Během trénování CNN jsou učeny váhy sítě (filtry) k detekci různých úrovní komplexnosti příznaků. Aby síť dokázala detekovat příznak nezávisle na jeho transformaci (pootočení, převrácení, posuvu), měla by být invariantní. Toho lze docílit tak, že neurony budou sdílet některé váhy (filtry).

Konvoluční operace je definována jako

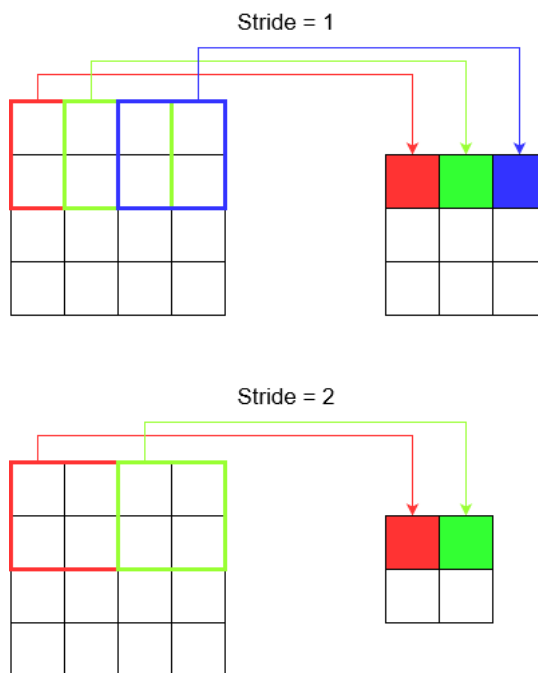
$$s[n] = (x \star w)[n] \quad (3.7)$$

kde  $x$  jsou vstupní data obvykle ve formě multidimenzionálního pole (tenzoru) a funkce  $w$  je konvoluční jádro. Výstupem operace je mapa příznaků. Konvoluční operaci lze definovat i pro vyšší dimenze. Máme-li například 2-D vstupní obrázek, pak i konvoluční jádro bude dvourozměrné. Přestože v mnoha knihovnách lze nalézt implementovanou operaci konvoluce, často se jedná o cross-korelaci.

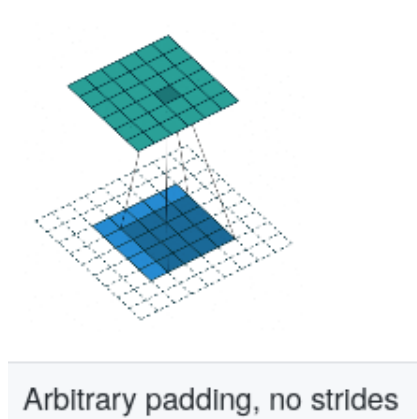
Průběh a chování konvoluce lze řídit hyperparametry *stride*, *padding* a *dilation*. Hodnota *stride* modifikuje velikost kroku konvolučního jádra. Výchozí hodnota kroku je 1. Chování při různých hodnotách lze vidět na obrázku 3.11.

Hodnota *padding* určuje, o kolik se mají rozšířit vstupní data. V případě malého rozměru konvolučního jádra může dojít k nechtěné změně rozměrů vstupních dat a proto se tato data po stranách rozšíří na takový rozměr, aby po konvoluci byl výsledný tvar v požadovaném rozměru. Vyplnění ukazuje obrázek 3.12.

**[[Tyto tri obrázky stylove sjednotit a dat je vedle sebe mozna]]**

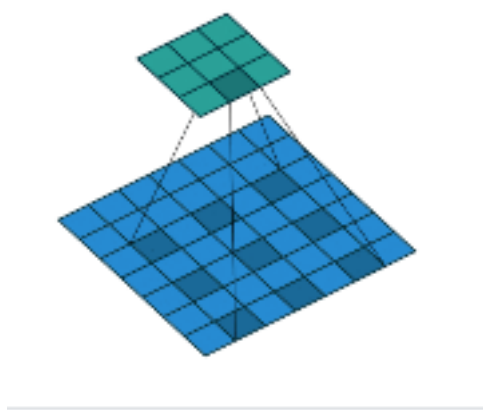


Obrázek 3.11: Hodnota stride ovlivňuje konvoluční krok jádra nad zpracovávanými daty. Čím větší hodnota stride, tím menší je výstupní rozměr dat



Obrázek 3.12: Nastavením hodnoty padding se zpracovávaná data na okrajích rozšíří o danou hodnotu a toto rozšíření bude vyplněno nulami či jakoukoli jinou hodnotou

Hodnota dilation podporuje rozšíření pole působnosti bez ztráty rozlišení nebo pokrytí dat konvolučním jádrem během konvoluční operace (viz obrázek 3.13). Zatímco pole působnosti roste exponenciálně, počet parametrů roste pouze lineárně. Těto vlastnosti se využívá například v temporal convolutional networks (TCN), kde dilation faktor nabývá hodnot  $1, 2, 4, 8, \dots$



Obrázek 3.13: Vizualizace hodnoty dilation upravující pole působnosti

### 3.4.1 Bottle-neck vrstva

[[Mozna, kdyz bude cas]]

## 3.5 Učení neuronových sítí

Cílem trénování je naučit neuronovou síť vykonávat nějaký úkol. Trénování založené na metodě učení s učitelem vyžaduje dostatečně velký dataset obsahující dvojice vstup–výstup a obnáší volbu objektivní funkce, kterou se během trénování snažíme minimalizovat.

Formálně je cílem učení aproximovat nějakou funkci  $f^*$ . Síti je předána vstupní hodnota  $x$ , pro kterou síť definuje mapování na výstupní hodnotu jako  $y = f(x; \theta)$ , kde  $\theta$  je parametr, který se síť učí tak, aby dosáhla nejlepší aproximace funkce [7, p. 163].

### 3.5.1 Objektivní funkce

Objektivní funkce (loss funkce, cost funkce) je funkce, která je během učení minimalizována nebo maximalizována, v závislosti na konkrétním úkolu a kontextu použití. Tato funkce mapuje událost či hodnoty na reálné číslo reprezentující [[spravit uvozovky]] "cenu", která je asociována s touto událostí či hodnotami. Je-li tato funkce použita pro optimalizační problém, pak je cílem hodnotu minimalizovat a funkci se pak říká loss funkce, protože její výstupní hodnota nám udává velikost chyby, která se počítá na základě rozdílu mezi výstupem sítě a odpovídajícími trénovacími daty z datasetu. Pak platí, že čím menší chyba, tím lépe síť provádí svůj úkol, na který je trénována. Použije-li se hodnota objektivní funkce s opačným znaménkem, tak lze funkci použít jako hodnotící metriku.

V rámci učení sítě je žádoucí, aby gradient objektivní funkce byl dostatečně velký (prudký) a předvídatelný. V takovém případě bude dobře sloužit pro účely trénování. V případě malého gradientu by funkce saturovala (byla by příliš plochá) a to by mohlo negativně ovlivnit trénování sítě [7].

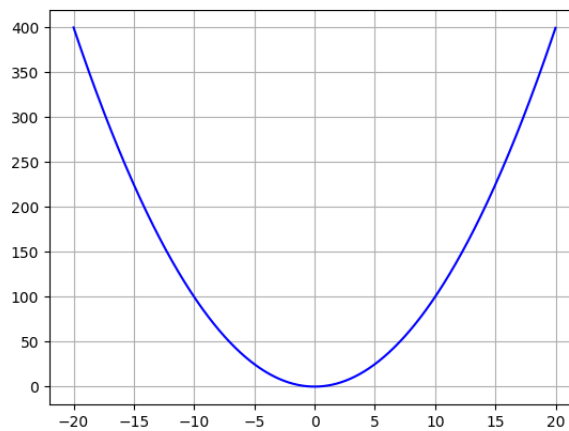
Nejčastěji se pro regresní problémy používají objektivní funkce *Mean Squared Error Loss* a *Mean Absolute Error Loss*.

## MSELoss (Mean Squared Error Loss)

MSELoss je nejčastěji používaná objektivní funkce při řešení regresních problémů. Vypočítá se jako

$$MSE = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{y}_i)^2 \quad (3.8)$$

kde  $N$  je počet trénovacích dat v datasetu,  $y_i$  je předpovězená hodnota a  $\hat{y}_i$  je odhadnutá hodnota. Vzhledem k druhé mocnině je výsledek vždy pozitivní nezávisle na znaménku hodnot  $y_i$  a  $\hat{y}_i$ , jak lze vidět na obrázku 3.14. Druhá mocnina také zajišťuje, že čím větší je rozdíl mezi předpovídanou a aktuální odhadnutou hodnotou, tím více se chyba projeví. Funkce dosahuje L2 regularizace.



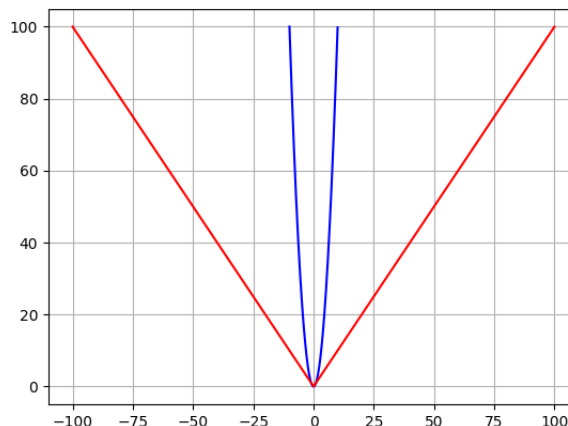
Obrázek 3.14: Graf objektivní funkce MSE Loss

## MAELoss (Mean Absolute Error Loss)

MAE je velmi podobná objektivní funkce jako MSE, ale s téměř opačnými vlastnostmi. Stejně jako MSE, ani tato funkce nenabývá negativní hodnoty, ale narozdíl od MSE, která má tuto vlastnost díky druhé mocnině rozdílu, MAE toho dosahuje tak, že rozdíl předpovídané hodnoty  $y_j$  a odhadnuté hodnoty  $\hat{y}_j$  je uzavřen v absolutní hodnotě. Vypočítá se jako

$$MAE = \frac{1}{N} \sum_{j=1}^N |y_j - \hat{y}_j| \quad (3.9)$$

Výhodnou MAE je její lineární průběh, takže se větší chyby projeví úměrně více než chyby menší, oproti MSE, kde je závislost kvadratická. Nevýhodou je, že kvůli absolutní povaze není diferencovatelná v hodnotě  $x = 0$ , což může mít negativní následky pro výpočet gradientu.



Obrázek 3.15: Graf objektivní funkce MAE Loss (červeně)

[[Hinge loss možná]]

### 3.5.2 Backpropagation

Backpropagation (backprop) je algoritmus typu učení s učitelem pro učení neuronových sítí. Algoritmus zpětně prochází neuronovou sítí a na základě hodnoty cost funkce  $J(\theta)$  počítá její gradient metodou *gradient descent*. Cílem je upravit všechny váhy v síti na základě toho, jak moc každá váha přispívá k celkové chybě. Algoritmus backpropagation počítá řetězové pravidlo (*chain rule*).

Algoritmus se skládá ze 3 opakujících se kroků. Prvním krokem je forward fáze, ve které se počítá pro každou dvojici vstupu a referenčního výstupu

[?] [?, p 197]

- zpětné šíření chyby - adaptační algoritmus, podíl neuronu na chybě, - 3 opakující se fáze učení:

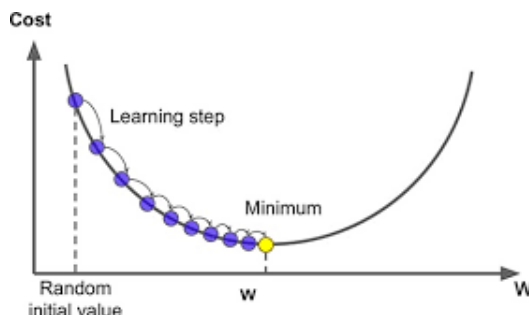
1) feedforward - dopředu 2) zpětné šíření chyby - Backpropagation - chain rule 3) úprava vah a biasu na základě chyby pomocí gradient descent

[[k poslednímu vodu zadání, ze diskutujte – tak lze jen napsat, že by mohlo toto toto, prostě ukázat, že jsem se nad tím minimálně zamyslel. mrknout na Tasnet studii, kde to řešili jak by to slo vyřešit, případně jestli bych našel nějaké studie které diskutují o zlepšení. Když bude čas, tak to pak i vyzkoušet a porovnat. toto Podsekce na konci experimentu.]]

### Gradient descent

Gradient descent je iterativní optimalizační algoritmus. Je určený pro hledání lokálního minima diferencovatelné objektivní funkce  $J(\theta)$ , kde  $\theta \in \mathbb{R}^d$  reprezentuje učené parametry modelu. Hledání lokálního minima probíhá úpravou parametrů modelu v opačném směru, než je hodnota gradientu. Výpočet je prováděn postupným pohybem ve směru největšího klesání, které je určeno zápornou hodnotou gradientu. Rychlost pohybu závisí na velikosti kroku, což udává *learning rate*. Správná volba *learning rate* ovlivňuje rychlost, jakou je nalezeno minimum funkce. Při nízké hodnotě bude výsledek přesnější, ale nalezení minima

bude výpočetně náročnější, protože v každém kroku se počítá nová hodnota gradientu. Při větší hodnotě je riziko, že minimum bude přeskočeno [18].



Obrázek 3.16: Ukázka algoritmu gradient descent při hledání minima dané funkce postupným přibližováním dle nastavené velikosti kroku

[[možna sem hodit tu rovnici updatu vah]] [[MOŽNA popsat více vliv LR na loss]]

### 3.5.3 Underfitting a generalizace

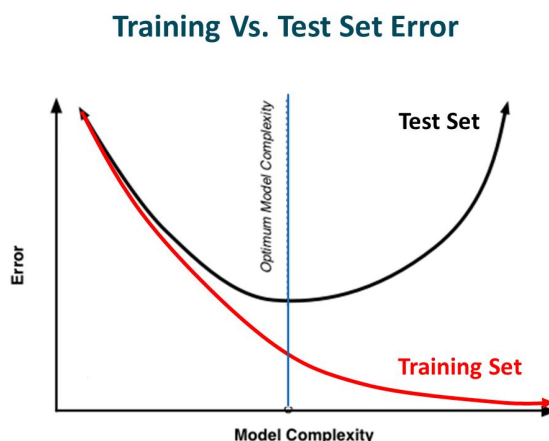
Generalizace, neboli *test error*, je vlastnost modelu, která udává, jak dobře model pracuje s dříve neviděnými daty, tedy s daty, které nebyly použity pro jeho učení. Během učení vzniká učicí chyba, která je trénováním postupně minimalizována, ale kvalita modelu je pak měřena na testovacích datech. Testovací chyba by ideálně měla být co nejbližší chybě na trénovací sadě. Z rozdílu mezi těmito chybovými hodnotami se dají diagnostikovat problémy s naučením sítě – podučení (*underfitting*) a přeučení (*overfitting*).

*Underfitting* nastane v případě, že model již nedokáže v důsledku jeho kapacity zmenšit dostatečně jeho chybovou hodnotu na trénovacím datasetu.

*Overfitting* nastane, když rozdíl mezi trénovací chybou a testovací chybou je příliš velký.

Kapacita je vlastnost modelu určující, kolik se toho dokáže model naučit. V případě nízké kapacity se model nedokáže naučit všechny příznaky z trénovacích dat a v případě vysoké kapacity se mohou přeučit zapamatováním si jejich příznaků a tím sice dosáhnou nízké trénovací chyby, ale špatných výsledků během testování [7, p107].



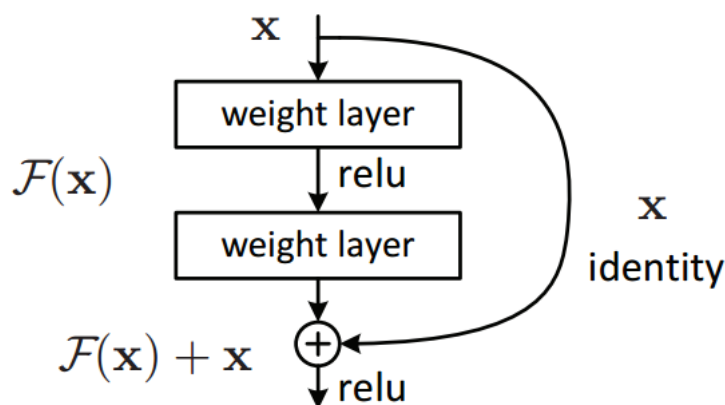


Obrázek 3.17: Vliv kapacity modelu trénovací a testovací chybu. Obrázek inspirován podle [3]

Regularizace je souhrn postupů, které říkají, jak modifikovat učící algoritmus tak, aby se zredukovala chyba generalizace, zatímco testovací chyba zůstane stejná [7, p117].

### 3.5.4 Reziduální spojení

Velmi hluboké neuronové sítě přinesly v mnoha aplikacích velmi dobré výsledky. S přibývajícím hloubkou se ale stává její trénování složitější. Jeden z problémů, který se při učení projevoval, byl explodující a mizející gradient [2]. Tento problém byl adresován zavedením normalizačních vrstev [8] a počáteční optimalizované inicializace, což dovolilo sítím konvergovat s využitím *stochastic gradient descent* se zpětnou propagací (*back-propagation*) [11]. Poté se objevil problém degradace, který nebyl zapříčiněn přetrénováním, a způsobil, že s přibývajícím hloubkou sítě její přesnost náhle prudce klesla, což indikovalo, že systémy nelze optimalizovat stejným způsobem a přidání dalších vrstev pouze zvýšilo trénovací chybu. Problém degradace lze řešit zakomponováním reziduálních spojení, které je zobrazeno na obrázku 3.18 [5].



Obrázek 3.18: Reziduální spojení mezi vstupem a výstupem stavebního bloku sítě [5]

**[[prepsat F na to specialni]]** Reziduální spojení, neboli identitní mapování, je zkratka mezi jednou či více vrstvami. Problém mizejícího gradientu řeší tím, že používá aktivace z předchozí vrstvy, dokud váhy následující vrstvy nejsou naučeny. Spojení lze formulovat jako  $F(x) + x$ , kde  $F(x)$  je výsledná hodnota transformace jednotlivými vrstvami a  $x$  je původní vstup, který vrstvy přeskočí a následně je sečten s výstupem  $F(x)$ . Výhodou je, že jeho použitím se nezvyšuje počet parametrů ani výpočetní náročnost [5][24].

### 3.5.5 Význam validační množiny v trénování

**[[Mozna ne jako subsekcí ale jen jako jeden odstavec]]** Většina algoritmů strojového učení má nějakou sadu hyperparametrů, kterou je upravováno chování algoritmu. Hodnoty hyperparametrů obvykle bývají nastavovány ručně ještě před spuštěním procesu učení a hodnota se v průběhu nemění, protože hodnoty by bylo obtížné optimalizovat. Některá nastavení se nicméně mohou stát hyperparametrem a být upravována během trénování, ale není vhodné je měnit na základě výsledku učení na trénovací sadě, protože by mohlo dojít k přetrénování (overfitting) v důsledku **[[CEHO??]]**. Pro tento případ potřebujeme validační sadu, která je odlišná od trénovací sady. Po každém zpracování trénovací sady následuje validační sada, po jejímž skončení jsou optimalizovány hyperparametry []. **[[[kniha 117-118] kap 5.3 = Hyperparameters and validation set]]** **[[najít jeste nejakej zdroj s popisem a pripadne nejaky zajimavejsi info.]]**

V této kapitole byly popsány základní principy neuronových sítí, jejich základní struktura a stavební bloky. V návaznosti na to byly zmíněny konvoluční sítě, které jsou pro tuto práci stěžejní. Dále byl popsán proces trénování sítí a algoritmy back-propagation a gradient descent. Nakonec byly zmíněny některé oblasti, které se zabývají optimalizacemi sítě, kladoucí si za cíl zlepšit například dobu trénování nebo kvalitu naučení sítě, což může mít za následek lepší výkon sítě při její evaluaci nebo umožnit nasazení ve zdrojově omezeném prostředí.

## Kapitola 4

# TasNet - Time-Domain Audio Separation Network

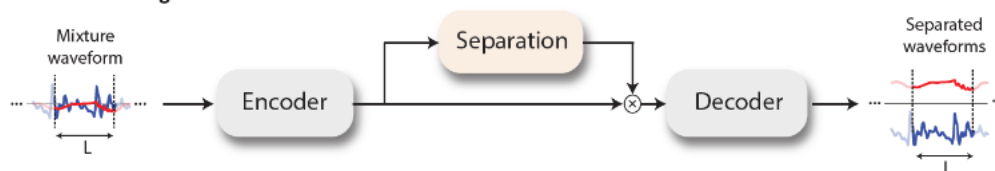
Tato kapitola vychází z referenční studie **TasNet: Surpassing Ideal Time-Frequency Masking for Speech Separation** [12] a popisuje její navrhnuté řešení pro separaci mluvčích v časové doméně.

Přestože metody pro zpracování řeči v takovém akustickém prostředí, ve kterém se současně prolíná mnoho řečových signálů, v poslední době zaznamenaly velké zlepšení, stále trpí mnoha nedostatky. Přesnost systémů, odezva a nároky na výpočetní výkon jsou u těchto metod natolik zásadní, že znemožňují nebo velmi omezují jejich nasazení mimo výzkumné prostředí, například v aplikacích, které by mohly zpracovávat řeč v reálném čase.

Většina dosavadních postupů byla založena na převodu směsi řečových signálů do časově-frekvenční (T-F) reprezentace (spektrogramu) pomocí short-time fourier transformation (STFT) [23]. Tento převod ale měl pro využití v reálném čase příliš vysokou odezvu a navíc T-F reprezentace nebyla optimální pro separaci mluvčích.

Pro překonání nedostatků předešlých metod byla navržena architektura fully-convolutional time-domain audio separation network (Conv-TasNet), založena na hlubokém učení a konvolučních neuronových sítích. Model prvně použije konvoluční enkodér k převodu krátkých segmentů směsi mluvčích na odpovídající nezápornou reprezentaci, která je optimalizovaná pro extrakci jednotlivých mluvčích. Samotné separace je docíleno aplikací masek na danou reprezentaci. Masky pro každého mluvčího pro každý segment v každém časovém kroku jsou odhadnuty v temporal convolutional network (TCN) [\[citace\]](#), která je tvořena opakující se posloupností konvolučních bloků se zvyšující se časovou dilatací. Po aplikaci masek jsou separovaní mluvčí rekonstruováni lineárním dekodérem. Tuto posloupnost operací zobrazuje zjednodušený obrázek 4.1. Dále v kapitole budou jednotlivé části popsány detailněji.

A. TasNet block diagram



Obrázek 4.1: Zjednodušený model architektury TasNet

Problém separace mluvčích na jednonanálové nahrávce lze definovat jako odhad  $C$  zdrojů mluvčích  $s_1(t), \dots, s_c(t) \in \mathbb{R}^{1 \times T}$  na diskrétním signálu směsi  $x(t) \in \mathbb{R}^{1 \times T}$ , kde  $T$  je délka nahrávky a kde

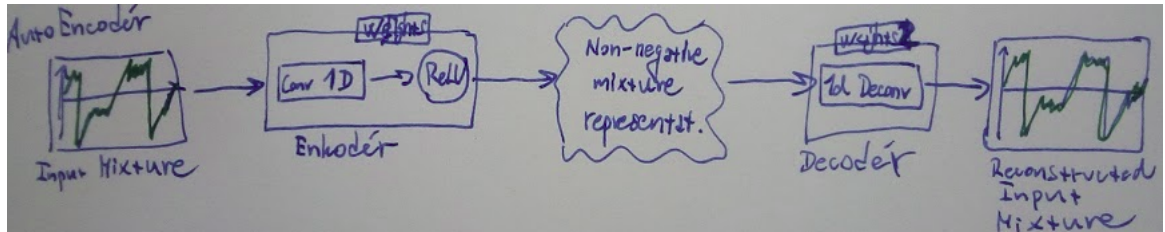
$$x(t) = \sum_{i=1}^C s_i(t) \quad (4.1)$$

Cílem je odhadnout  $s_i(t), i = 1, \dots, C$  ze signálu směsi  $x(t)$ .

[[sem možná tabulka hyperparametrů a jejich pojmenování **X, R, N, L, P, ...**]]

## 4.1 Konvoluční auto-enzodér

Konvoluční auto-enzodér převádí vstupní segmenty nahrávky na nezápornou reprezentaci a následně zase zpět z reprezentace na původní nahrávku, jak ukazuje obrázek 4.3.



Obrázek 4.2: Schéma konvolučního autoenzodéru

### 4.1.1 Enkódování směsi

[[spravit znak pro konvolucni operator v math env]] Každý segment směsi o délce  $L$ ,  $x_k \in \mathbb{R}^{1 \times L}$ , kde  $k = 1, \dots, \frac{T}{L}$ , je transformován na nezápornou reprezentaci  $w \in \mathbb{R}^{1 \times N}$  pomocí 1-D konvoluční operace jako

$$w = \text{ReLU}(x \star U) \quad (4.2)$$

kde  $U \in \mathbb{R}^{N \times L}$  obsahuje  $N$  vektorů, každý délky  $L$ , které reprezentují bázevské funkce enkodéru. Operace  $\star$  značí konvoluční operaci. ReLU je nelineární aktivační funkce, která byla blíže popsána v podkapitole 3.3.

### 4.1.2 Dekódování extrahovaných mluvčích

Pro převod z reprezentace zpět do podoby audio nahrávky slouží lineární dekodér. Pomocí 1-D dekonvoluční operace rekonstruuje původní signál  $x$  jako  $x \in \mathbb{R}^{1 \times L}$ . Tuto operaci lze definovat jako

$$\hat{x} = wV \quad (4.3)$$

kde každý řádek v matici  $V \in \mathbb{R}^{N \times L}$  představuje jednu bázevou funkci dekodéru s délkou  $L$ .

## 4.2 Separační modul

Cílem separačního modulu je najít váhovanou funkci, neboli masku, pro každého zdrojového mluvčího pro každý výstup enkodéru v každém časovém kroku. Formálně lze zapsat, že cílem

je odhadnout  $C$  masek  $m_i \in \mathbb{R}^{1 \times N}, i = 1, \dots, C$ , kde  $C$  představuje počet mluvčích ve směsi. Vektory masek  $m_i$  mají takové omezení, že  $\sum_{i=1}^C m_i = 1$ , kde  $1$  je jednotkový vektor v  $\mathbb{R}^{1 \times N}$ . Toto omezení garantuje, že rekonstruované zdroje po sečtení zformují původní směs  $\hat{x} = \sum_{i=1}^C \hat{s}_i$ .

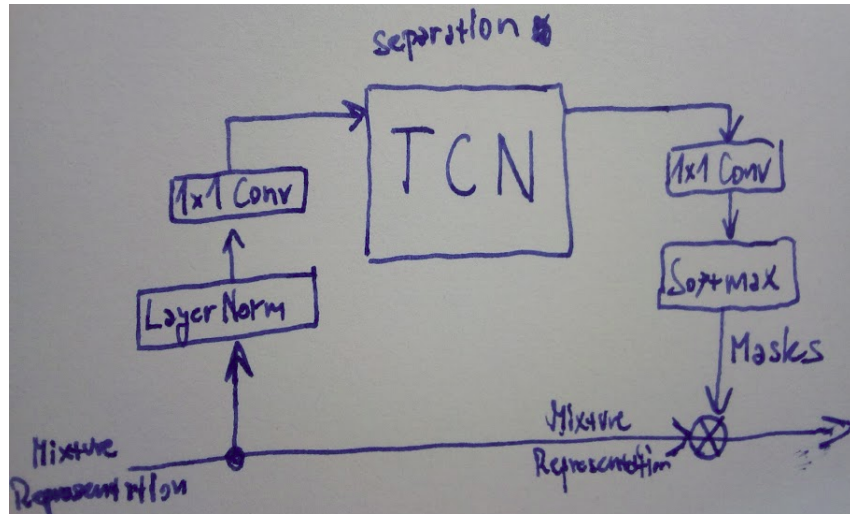
Separace je pro každý zdroj provedena vynásobením odpovídající masky  $m_i$  s nezápornou reprezentací (výstupem enkodéru)  $w$ , jako

$$d_i = w \odot m_i \quad (4.4)$$

kde  $d_i \in \mathbb{R}^{1 \times N}$  je reprezentace každého ze zdrojů a operace  $\odot$  je **[[not sure]]** vektorový součin.

Nakonec jsou preprezentace  $d_i$  rekonstruovány zpět na korespondující nahrávky zdrojů  $\hat{s}_i, i = 1, \dots, C$  pomocí dekodéru jako

$$\hat{s}_i = d_i V \quad (4.5)$$



Obrázek 4.3: Schéma kompletního separačního modulu, jehož vstupem je nezáporná reprezentace směsi mluvčích

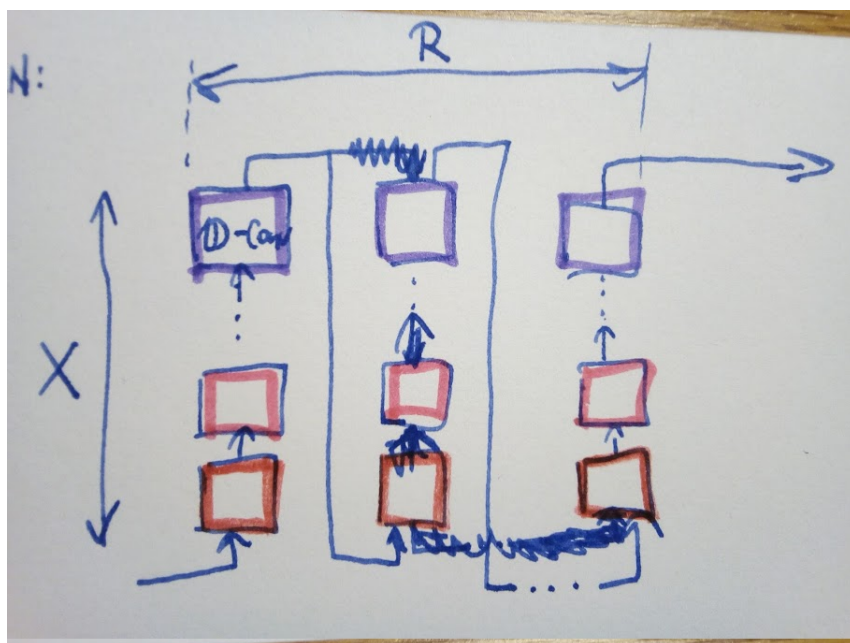
**[[obrazky mozna nektere zcuknout k sobe nebo vedle sebe]]**

Na začátku separačního modulu je přidána lineární operace  $1 \times 1$ -conv jako bottleneck, která určuje počet kanálů na vstupu a výstupu následující sekvence konvolučních bloků.

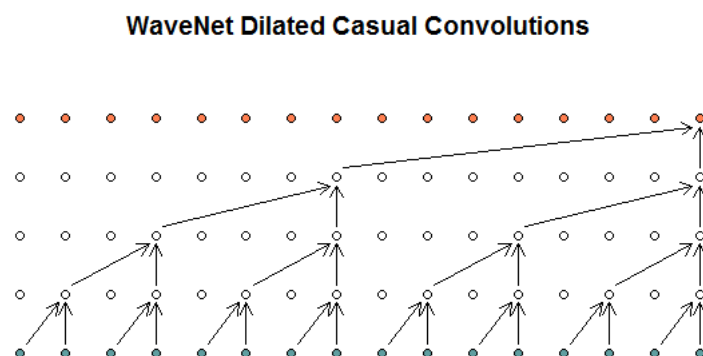
#### 4.2.1 Temporal Convolutional Network

Jádrem separační části je *temporal convolutional network* skládající se z naskládaných 1-D konvolučních bloků s časovou dilatací jak lze vidět na obrázku 4.4. Každá vrstva v TCN obsahuje bloky se zvyšující se časovou dilatací, která se zvyšuje exponenciálně v závislosti na parametru  $X$ , který udává počet konvolučních bloků **[[v jedné vrstvě]]** a nabývá hodnot  $1, 2, 4, \dots, 2^{X-1}$ . Taková sekvence bloků je opakována  $R$ -krát. Exponenciální růst dilatačního faktoru zajišťuje dostatečně velké okno pro využití výhod dlouhých časových závislostí v signálu řeči. Obrázek 4.5 zobrazuje rostoucí časovou dilataci na vstupních vzorcích signálu.

Výstup posledního bloku posledního opakování v TCN je předán  $1 \times 1$  konvoluční vrstvě, která má  $N \times C$  filtrů a nakonec aktivační funkci softmax k odhadu  $C$  vektorů masek pro každého z  $C$  cílových mluvčích.



Obrázek 4.4: Jádru separačního modulu – naskládané konvoluční bloky s časovou dilatací odhadující masky na základě nezáporné reprezentace[[**pridat dilation factor 2 mocnina X minus 1 do obrazku stacked-resblocks jpg**]]



Obrázek 4.5: Vizualizace časové dilatace

#### 4.2.2 Konvoluční bloky - možná ne jako podkapitolu

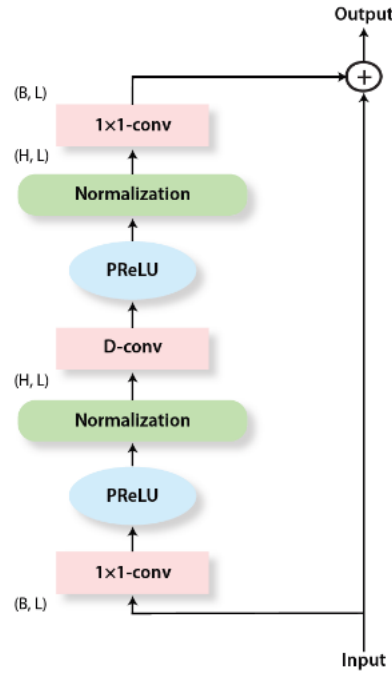
**[[Pozn: vrstvy bloku a residuální spojení]]** Navržená architektura dále nahradila standardní konvoluci uvnitř  $1 \times 1$  konvolučních bloků za *depthwise separable convolution* (S-conv), která pomáhá snížit počet parametrů a ukázala se jako efektivní ve zpracování obrazu **[[citace 26 27]]**.

*Depthwise separable convolution* se skládá ze dvou, po sobě jdoucích, operací – depthwise convolution (D-conv) a standardní konvolucí  $1 \times 1$ -conv s velikostí konvolučního jádra 1:

$$D - conv(Y, K) = \text{concat}(y_j \text{conv} k_j), j = 1, \dots, N \quad (4.6)$$

$$S - conv(Y, K, L) = D - conv(Y, K) \text{conv} L \quad (4.7)$$

kde  $Y \in \mathbb{R}^{G \times M}$  je vstup do S-conv,  $K \in \mathbb{R}^{G \times P}$  je konvoluční jádro o velikosti  $P$ , dále  $y_j \in \mathbb{R}^{1 \times M}$  a  $k_j \in \mathbb{R}^{1 \times P}$  jsou řádky matic  $Y$  a  $K$ .  $L \in \mathbb{R}^{G \times H \times 1}$  je konvoluční jádro o velikosti 1. Operace  $1 \times 1$ -conv se chová jako plně propojená vrstva a transformuje příznaky do potřebných rozměrů.



Obrázek 4.6: Jeden konvoluční blok obsahující reziduální spojení mezi vstupem a výstupem **[[graficky zde ukázat a zakroužkovat co je sconv a co je dconv]]**

V každém konvolučním bloku jsou operace  $1 \times 1$ -conv a D-conv následovány nelineární aktivační funkcí PReLU [6], která byla popsána v kapitole 3, a normalizační vrstvou. Konvoluční bloky dále obsahují reziduální spojení [5] mezi vstupem a výstupem, které výrazně zjednodušuje trénování velmi hlubokých neuronových sítí díky ponechané referenci na vstupní netransformovaná data.

**[[zvýraznit na obrázku konvolučního bloku residuální spojení a hodit tam sipku nebo něco ze to je to spojení]]**

**[[Kam uvozovky, kam bold, kam italic ZEPTAT SE]]**

### 4.3 Normalizace

Při trénování neuronových sítí může docházet k fenoménu jménem **[[uvozovky]]** vnitřní kovarianční posuv *internal covariate shift*, kvůli kterému je nutné pečlivě inicializovat parametry a volit menší hodnotu learning rate, což zpomaluje trénování sítě. K adresování problému slouží zakomponování normalizace do architektury neuronové sítě. Normalizace je aplikována na každý mini-batch trénovacích dat **[8][1]**.

Konvoluční bloky obsahují normalizační vrstvy, které mohou významně ovlivnit výkon sítě. Ve studii bylo experimentováno s channel-wise layer normalization (cLN) **[[cite]]**, global layer normalization (gLN) a batch normalization (BN) **[8]**.

Metoda cLN je aplikována na vstup separačního modulu pro zajištění invariance při změně měřítka vstupních dat. Tato metoda je vhodná pro použití v kauzální i nekauzální konfiguraci, je-li použita v konvolučních blocích. Normalizace cLN je aplikována na každý segment  $y_k$ .

Metoda gLN je aplikována globálně na každý příznak na rozměry kanálu i času. Tuto normalizaci lze použít pouze při nekauzální konfiguraci, protože výpočet probíhá na základě celého vstupu. Při kauzální konfiguraci může být použita také metoda BN.

Normalizace BN je počítána následovně:

$$y = \frac{x - E[x]}{Var[x] + \varepsilon} * \gamma + \beta \quad (4.8)$$

$$E[x] = \frac{1}{NT} \sum_{NT} Y \quad (4.9)$$

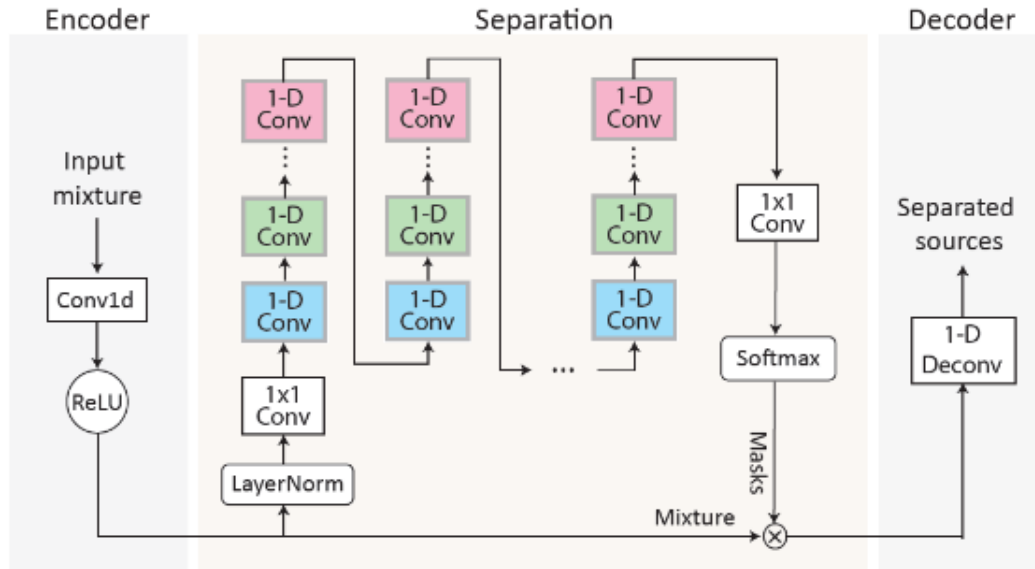
$$Var(Y) = \frac{1}{NT} \sum_{NT} (Y - E[Y])^2 \quad (4.10)$$

**[[Asi jeste popsat jednotlivé promenne ve vzorcich]]**

Jak vypadá kompletní schéma architektury po sestavení a propojení jednotlivých modulů lze vidět na obrázku **4.7**.



## B. System flowchart



Obrázek 4.7: Schéma architektury TasNet

V této kapitole byla popsána navržená architektura Conv–TasNet pro separaci navzájem se překrývajících mluvcích v časové doméně, skládající se z konvolučního auto–enkodéru, který převádí vstupní směs mluvcích na nezápornou reprezentaci, ze které je následně v separačním modulu složeném z konvolučních bloků se zvyšující se časovou dilatací odhadnuty multiplikativní funkce (tzv. masky) pro každý ze zdrojů. Ty jsou nakonec aplikovány na reprezentaci a výsledek lineárním dekodérem převeden zpět na hlasovou, nyní již separovanou nahrávku.

## Kapitola 5

# Implementace sítě

Pro implementaci neuronových sítí vzniklo mnoho frameworků jako PyTorch, Tensorflow, Keras a další, které umožňují síť poskládat z již předdefinovaných modulů. Jejich vlastnosti lze modifikovat dle potřeby pomocí argumentů při jejich instanciaci nebo při jejich použití. Kromě těchto bloků a mnoha dalších možností frameworky obsahují i metody pro práci s daty, funkce pro vykreslování grafů a pro matematické operace.

Mým úkolem bylo implementovat neuronovou síť podle architektury TasNet [12] pro separaci mluvčích v časové doméně. V rámci zadání jsem si zvolil použít framework s názvem PyTorch<sup>1</sup>, který má kvalitní dokumentaci a aktivní uživatelskou základnu. PyTorch funguje nad jazykem Python<sup>2</sup> a v celé práci používám jeho aktuálně nejnovější verzi (Python 3.8.2). Mimo jiné podporuje práci s daty, včetně implementace vlastního nahrávání a manipulace s daty. Základní jednotkou, se kterou síť pracuje, je tensor. Ve frameworku PyTorch je **torch.Tensor** třída reprezentující multidimenzionální matici obsahující data jednoho typu. Na tensors jsou převáděna data během načítání datasetu a následně předávána síti ke zpracování. Samotná implementace probíhala na systému Kali linux (4.19.0-kali5-amd64), který je založený na systému Debian. Aplikace byla vyvíjena v textovém editoru Vim<sup>3</sup> a byla během implementace spouštěna a testována v příkazovém procesoru Bash<sup>4</sup> (Bourne Again SHell) verze 5.0.16. Bash jsem zvolil i pro implementaci některých pomocných skriptů kvůli jednoduchému spouštění z terminálu a možnostem automatizace některých opakujících se příkazů. Pro správu verzí jsem používal verzovací systém Git<sup>5</sup>.

Trénování sítě probíhalo z počátku na CPU *intel i5* mého osobního notebooku *Lenovo Y50-70*, což se ale ukázalo jako naprosto nevhodné vzhledem k náročnosti výpočtu kvůli nedostatečnému výkonu. Začal jsem tedy používat službu *Google colab*, která poskytuje na omezenou dobu (12 hodin) stroje, které navíc obsahují technologii *cuda*, což mi umožnilo výpočet provádět na GPU, čímž se trénování zrychlilo několikanásobně. Ještě více pomohlo předplacení služby, čímž se mi prodloužila doba, po kterou jsem mohl síť trénovat, na 24 hodin a zvýšila se šance na získání silnějších strojů. Modely byly trénovány na GPU *T80* a *T100* v závislosti na tom, která z nich mi zrovna byla přidělena pro dané sezení (*session*).

---

<sup>1</sup><https://pytorch.org/>

<sup>2</sup><https://www.python.org/>

<sup>3</sup><https://www.vim.org/>

<sup>4</sup><https://www.gnu.org/software/bash/>

<sup>5</sup><https://git-scm.com/>

Pro vyhodnocování výkonu sítě byly použity metriky SDR, PESQ a STOI. Pro výpočet SDR jsem použil zdrojový kód nacházející se na githubu <sup>6</sup>. Metriky PESQ <sup>7</sup> a STOI <sup>8</sup> byly nainstalovány pomocí balíkového nástroje pip <sup>9</sup>.

**[[zamyslet se co tu chci ukazat v tom algu, pripadne vypustit parsovani, in-**

```

parsování argumentů programu;
instanciace modelu TasNet;
if args.checkpoint! = NULL then
    | načtení checkpointu;
end
instanciace třídy SegmentDataset a dataloaderu pro trénovací data;
instanciace třídy SegmentDataset a dataloaderu pro validační data;
for epocha v args.Epochs do
    /* trénování */
1   loss = 0;
2   for směs_mluvčích z trénovacího datasetu do
    | /* cyklus načítá z dataloaderu segmenty nahrávek, které jsou
    |    generovány v třídě SegmentDataset */
3   | odhadovan_separace = TasNet(sms_mluvch);
instanciaci atd]] 4   | loss = sisnr(optimln_separace, odhadovan_separace);
5   | propagace chyby a úprava vah;
6   end
    /* validace */
7   vypnout gradient;
8   for směs_mluvčích z validačního datasetu do
9   | odhadovan_separace = TasNet(sms_mluvch);
10  | validacni_loss = sisnr(optimln_separace, odhadovan_separace);
11  | if hodnota validační loss neklesla potřetí v řadě then
12  | | learning_rate = learning_rate/2;
13  | end
14  end
15  na konci epochy uložit checkpoint;
end

```

**Algoritmus 1:** Zjednodušený algoritmus běhu programu pro trénování sítě

## 5.1 Implementace tříd a modelu

Model neuronové sítě byl logicky rozdělen na několik dílčích částí, z nichž každá je reprezentována jednou třídou. Jednotlivé třídy obsahují metody pro nahrávání a transformaci dat, segmentaci nahrávek a další jsou přímo stavebním blokem neuronové sítě. Některé parametry předané třídám při jejich instanciaci jsou hyperparametry, jejichž hodnoty byly určeny parametry z příkazové řádky při volání programu. Každá z tříd bude popsána dále v textu.

<sup>6</sup>[https://github.com/craffel/mir\\_eval](https://github.com/craffel/mir_eval)

<sup>7</sup><https://github.com/ludlows/python-pesq>

<sup>8</sup><https://github.com/mpariente/pystoi>

<sup>9</sup><https://pypi.org/project/pip/>

## Třídy SegmentDataset a AudioDataset

Třídy **SegmentDataset** a **AudioDataset** slouží pro načítání korespondujících trojic nahrávek směr–zdroj1–zdroj2 z adresáře, kde je uložen dataset. Nahrávky jsou ve formátu wav a vzorkovány na 8KHz. Třídy tato data dále transformují na tenzory a normalizují na jednotkovou varianci a zero–mean.

Tyto třídy používají některé pomocné funkce definované v souboru **tools.py** pro dorovnání nedostatečně velkých nahrávek ve zpracovávaném mini–batchi nulami na velikost nejdelší z nich, protože data obsažená v mini–batchi musejí mít stejnou velikost.

Při spuštění programu jsou nahrány názvy wav souborů do polí, které jsou poté buď náhodně či sekvenčně procházeny a postupně dochází k načítání trojic nahrávek na aktuálním indexu.

Obě třídy jsou periodicky volány instancí třídy **Dataloader**<sup>10</sup>, která v cyklu posílá požadavky na data ke zpracování. Rozdíl mezi třídami je ten, že třída **AudioDataset** předává trojice nahrávek v původním stavu a délka epochy tak trvá tolik iterací, kolik je nahrávek ve zpracovávaném datasetu a je kvůli tomu používána během testování sítě, kde chceme vyhodnocovat celé nahrávky. Třída **SegmentDataset** je použita pro trénovací a validační dataset, kde jsou nahrávky rozdělovány na segmenty o parametrem dané délce a až poté jsou poskytnuty ven. Segmentace je detailně popsána v podkapitole 5.2.

Po vyčerpání datasetu obě třídy vyvolají vyjímku **StopIteration**, která je detekována **Dataloaderem** a ten ukončí cyklus, ve kterém se dataset zpracovával. Tím v případě trénování končí učení sítě a začíná validace, v případě validace končí jedna epocha a v případě testování končí vyhodnocování sítě.

## Třída Tasnet

Třída **Tasnet.py** reprezentuje model neuronové sítě. Atributy třídy představují jednotlivé vrstvy neuronové sítě. Sít obsahuje konvoluční vrstvy, enkodér, dekodér a *temporal convolutional network*, která je tvořena sekvencí konvolučních bloků, které jsou rozebrány později. Konvolučním vrstvám jsou nastaveny hodnoty *padding* a *stride*, které ovlivňují chování konvolučních vrstev. V této třídě jsou také inicializovány váhy (filtry) konvolučních vrstev algoritmem Xavier [4]. Ve funkci *forward()*, která je volána, když jsou instanci sítě předána data ke zpracování, jsou sekvenčně volány jednotlivé vrstvy. Během zpracování předaných dat je v této sekvenci vypočítána maska a aplikována na zpracovávanou směs. Tím vzniknou dvě separované nahrávky, které jsou předány na výstup v podobě tensoru.

## Třída ResBlock

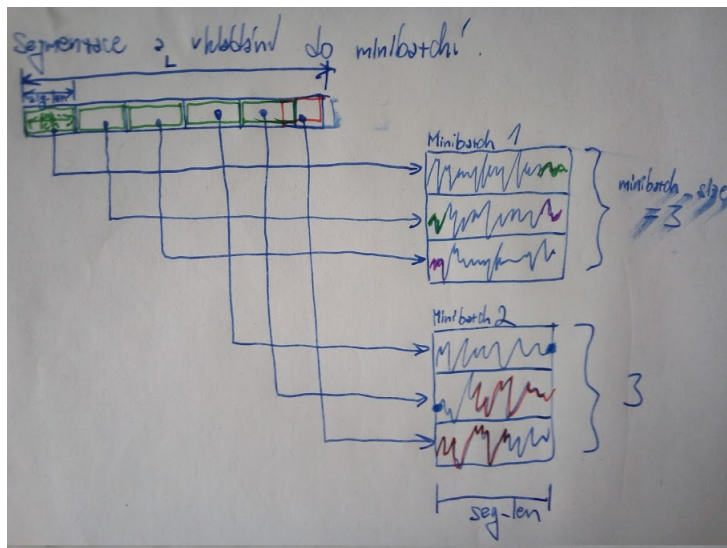
Třída **ResBlock** reprezentuje jeden konvoluční blok, který je opakován se zvyšující se dilatací v *temporal convolutional network* v separační části určené k odhadu masek. Podobně jako třída **Tasnet**, i tato třída obsahuje inicializaci a zřetězení vrstev. V inicializaci jsou nastaveny hodnoty pro konvoluční operace jako počet konvolučních jader, časová dilatace a počet kanálů. Tyto operace jsou ve funkci *forward* zřetězeny do sekvence konvolučních operací a aktivačních funkcí, podle vzorové architektury **TasNet**. Třída **ResBlock** je specifická svým residuálním spojením, kde dochází k rozdělení dat při zpracování, kdy jedna kopie je transformována operacemi v bloku a na konci sečtena s daty, které transformacemi neprošly a tento součet je výstupem bloku. Residuální spojení jsou podrobněji popsána v kapitole 3.

<sup>10</sup><https://pytorch.org/docs/stable/data.html#torch.utils.data.DataLoader>

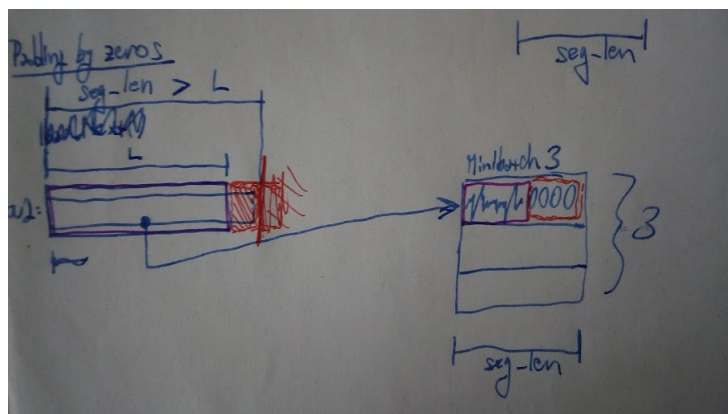
## 5.2 Segmentace nahrávek

Pro účely trénování sítě byly vstupní nahrávky rozdělovány na segmenty o délce 4, případně 2 sekund. Délku segmentů  $L_s$  lze nastavit přes argument `--segment_length` při volání skriptu *train.py*. Tato hodnota reprezentuje jeden z hyper-parametrů sítě. Výchozí hodnota je 4 sekundy. Při nižších hodnotách se prodlužuje délka trénování, protože se z jedné vstupní nahrávky o délce  $L$  vygeneruje až  $L/L_s$  segmentů a tím narůstá počet dat, který síť musí zpracovat.

Poslední segment by byl kratší, pokud by délka nahrávky byla kratší než délka segmentu. V takovém případě je segment doplněn nulami do délky ostatních segmentů jak lze vidět na obrázku 5.2. Toto doplnění je nutné kvůli dávkovému zpracování (**batches**), které neumožňuje rozdílnou délku elementů v dávce. Pokud je délka nahrávky delší než délka segmentu, tak se poslední segment vezme od konce bez nutnosti ho pak nulami doplňovat. Skládání segmentů do minibatche zobrazuje obrázek 5.1. Při vzorkovací frekvenci  $f_s = 8000\text{Hz}$ , obsahují segmenty  $f_s * L_s$  vzorků vstupní nahrávky.



Obrázek 5.1: Dělení nahrávky na segmenty o délce `segment_len` a jejich následné vkládání do minibatche. Lze si všimnout, že pokud poslední segment není dostatečně dlouhý, tak se vezme od konce nahrávky a jeho počáteční vzorky budou duplicitní s posledními vzorky předešlého segmentu



Obrázek 5.2: Je-li zpracovávána nahrávka, jejíž délka je kratší než je délka segmentu, je tento segment doplněn z prava nulami do požadované délky

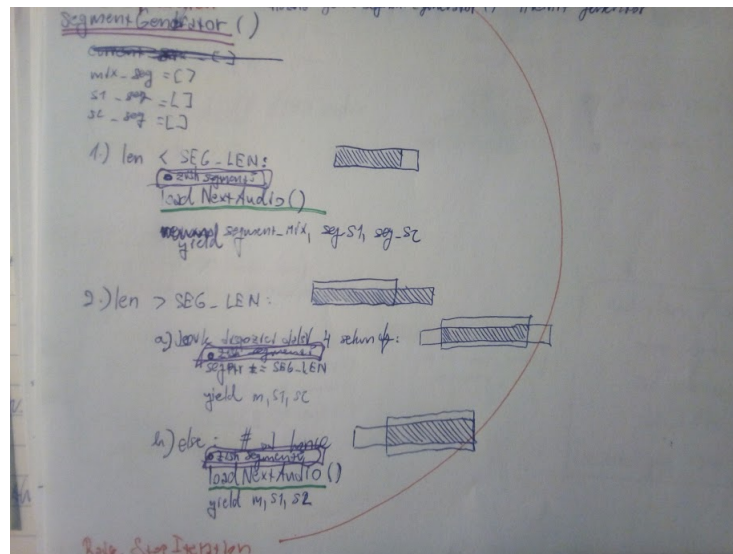
**[[možna tento odstavec jeste nejak upravit a doplnit obrazek kde je nactani asi 3 audio nahravek]]** Algoritmus 2 a obrázek 5.3 zobrazují proces segmentace. Funkce `segmentGenerator()`, která slouží jako generátor, pracuje současně se třemi nahrávkami – se směsí dvou mluvčích, s nahrávkou prvního mluvčího a s nahrávkou druhého mluvčího, které jsou do atributů `current_mixture`, `current_source1` a `current_source2` třídy `SegmentDataset` nahrávány funkcí `loadNextAudio()` poté, co z předchozí trojice nahrávek již nelze vygenerovat další segmenty.

```

Function segmentGenerator() : int[] is
    s1_segment = [];
    segptr = 0;
    while je dostupná další trojice nahrávek do
        if current_mixture_len < SEGMENT_LENGTH then
            /* aktuální nahrávka je kratší než délka segmentu */
1         s1_segment = current_source1[:];
2         yield s1_segment;
        else
            if segptr + SEGMENT_LENGTH < current_mixture_len then
                /* z aktuální nahrávky vzít segment */
3                 s1_segment = current_source1[(segptr+SEGMENT_LENGTH)];
4                 segptr += SEGMENT_LENGTH;
5                 yield s1_segment;
            else
                /* nelze již načíst celý segment, takže se vezme od konce */
6                 start_index = (current_mixture_len - SEGMENT_LENGTH);
7                 s1_segment =
                    current_source1[start_index:current_mixture_len];
8                 yield s1_segment;
            end
        end
    end
end

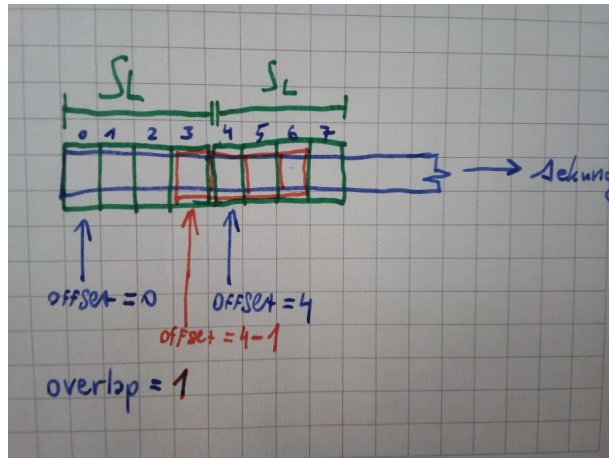
```

**Algoritmus 2:** Algoritmus segmenace nahrávek používá výhody Python generátoru, který, narozdíl od běžných funkcí, po vrácení hodnoty příkazem `yield` neztrácí svůj vnitřní stav a při jeho dalším zavolání pokračuje tam, kde skončil. V uvedeném příkladě je znázorněna pouze segmentace nahrávky jednoho mluvčího, ale analogicky se segmentují i nahrávky směsi mluvčích a nahrávky druhého mluvčího



Obrázek 5.3: Způsob segmentace nahrávek

Kód segmentace také umožňuje určit délku, o kterou se segmenty budou překrývat. Těto funkcionality je docíleno zavedením proměnné **overlap**, která obsahuje délku překrytí a při segmentaci je od proměnné **offset**, který určuje začátek následujícího segmentu, tato hodnota odečtena, jak ukazuje obrázek 5.4. Při trénování je překrytí nulové, protože by to zvýšilo počet trénovacích dat jejich částečnou duplikací, což by negativně ovlivnilo délku trénování.



Obrázek 5.4: Segmentace nahrávek s nenulovým překrytím

### 5.3 Pomocné skripty

Během implementace se vyskytla potřeba zjednodušit si některé repetitivní úkony a jelikož upřednostňuji práci z terminálu, tak jsem si pro ně naprogramoval několik skriptů za použití jazyků Python3 a Bash.



Skript **gupload.sh** pro odesílání zdrojových souborů na *Google Drive*<sup>11</sup>, ze kterého *Google Colab*<sup>12</sup> načítá data a kód, který lze následně po částech spouštět z grafického rozhraní. K interakci s Google Drive jsem použil nástroj *rclone*<sup>13</sup>, který nabízí mnoho operací pro vzdálenou práci s Google Drive z terminálu. S pomocí *rclone* lze stahovat či nahrávat data, vzdáleně měnit adresářovou strukturu, mazat, vytvářet a upravovat soubory a podobně.

Skript **gdownload.sh** je určený pro snadné stahování dat z Google Drive. Stejně jako předchozí skript, používá operace poskytované nástrojem *rsync*. Při zavolání stáhne adresář obsahující data z posledního trénování, nebo adresář, jehož název je předán jako argument skriptu v příkazové řádce. Stažený adresář obsahuje checkpointy, soubory obsahující výstup trénování, soubor obsahující dvojici hodnot  $[zpracovno_nahrvek, loss]$ , které lze jednoduše vykreslit do grafu, a nakonec soubor s výsledky testování obsahující hodnoty metrik pro testovaný checkpoint.

Dále jsem vytvořil několik jednoduchých skriptů pro rychlé spuštění trénování a testování modelu s různými přednastavenými parametry. Bylo tak možno mít mezi skripty rozdílné cesty k adresářům pro účely trénování a ladění chyb nebo třeba nastavení času, po kterém se proces má sám ukončit. Pro spuštění trénování modelu tak lze zavolat skript **nntrain.sh**, pro testování skript **nnctest.sh** a pro inferenci (separaci mluvčích na předané nahrávce) skript **nninference.sh**.

Implementace samotného modelu ve frameworku PyTorch nebyla náročná. Vzhledem k netriviálnímu úkolu jsem musel implementovat některé části, jako třeba segmentace nahrávek, sám. To přineslo mnoho potíží při odstraňování chyb v samotném modelu. Síť jsem prvně trénoval a ladil na CPU svého notebooku, takže probíhalo velmi pomalu (v řádu desítek hodin i při nízké velikosti modelu), a některé chyby se objevily až v průběhu trénování nebo dokonce na jeho konci. Situaci razantně zlepšil přechod na Google Colab, kde jsem síť trénoval na poskytnutých GPU a trénování tak zabralo **[["pouze"(opravit uvozovky)]]** několik hodin, takže se případná chyba objevila dříve, ale i přesto mi ladení sítě zabralo mnoho času. Pro další usnadnění samotné implementace jsem vytvořil několik skriptů, které mi usnadnily spouštění a ladení sítě.

---

<sup>11</sup><https://www.google.com/drive/>

<sup>12</sup><https://colab.research.google.com/notebooks/intro.ipynb>

<sup>13</sup><https://rclone.org/>

## Kapitola 6

# Experimenty a vyhodnocení

- trenovani s ruznymi hyperparametry, uspesnost a tabulky s hyper parametry a dosazenymi vysledky a hodnotami sisnr, sdr atd. - model size comparison. - porovnaní s výsledky ze studie - obrázky separovaných mluvčích - signalu. - spektra - grafy trenovani loss a výsledkuu. - pametova narocnost modelu

Neuronové sítě představují mocný nástroj mnoha využití, který ale musí být naučen, jak danou činnost provádět. Trénování sítí je výpočetně náročný úkol, jehož náročnost se, krom dalších parametrů, může lišit v závislosti na velikosti datasetu a sítě. Obecně čím větší dataset a rozměr sítě, tím déle trénování trvá. Je proto podstatné, na jakém stroji je síť učena a od toho se odvíjí čas trénování.

Po dostatečném natrénování probíhá vyhodnocení sítě. Výsledky trénování jsou měřeny metrikami, které se volí na základě problému, který se síť učila řešit. Výsledky testování sítě jsou následně porovnávány v závislosti na velikosti daného modelu. Velikost modelu je ovlivněna hyper-parametry, které určují například počet vrstev sítě, velikost segmentace a další vlastnosti.

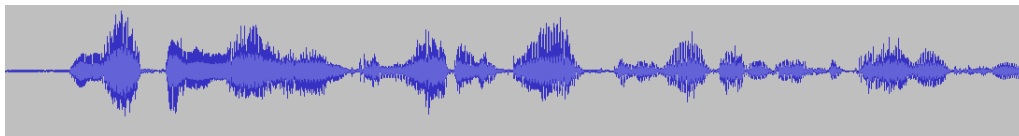
Cílem experimentů je zjistit, které z hyper-parametrů mají největší vliv na kvalitu separace a tedy na výsledky vyhodnocení, a jak moc lze model zmenšit, aby stále dával dostatečně dobré výsledky separace mluvčích. Tato kvalita je měřena primárně metrikami STOI, SI-SNR a PESQ, které byly popsány v kapitole 2. [\[\[ref na zdroje u metrik\]\]](#)

### 6.1 Dataset

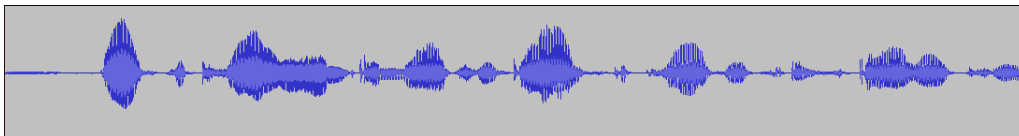
Dataset je množina dat, na kterých je síť učena a testována. Dataset se dá rozdělit na 3 podmnožiny, které jsou navzájem exkluzivní. Trénovací dataset je určen k trénování sítě a bývá největší. Validační dataset se prochází po zpracování trénovacího datasetu v průběhu trénování a slouží pro ověření, že hodnota loss na dosud neviděných datech stále klesá a nedochází tak k přetrénování. Při zpracování tohoto datasetu je síti zabráněno v učení. Poslední je testovací dataset, na kterém probíhá vyhodnocení sítě, tedy výpočet hodnot jednotlivých metrik.

Trénování a vyhodnocení modelu proběhlo na množině jednokanálových nahrávek směsí dvou mluvčích. Množina byla vygenerována náhodným výběrem různých mluvčích z Wall Street Journal (WSJ0) a vytvořením směsí mluvčích z těchto dvou náhodně vybraných nahrávek. Celková délka trénovacích dat je přes 10 hodin a přes 6 hodin validačních dat. Nahrávky jsou převzorkovány na 8kHz a během trénování zarovnány na zero means a jednot-

kovou varianci [[[studie str 5 Dataset]][49 - ze studie odkaz na script na generování a popis na netu]]].



Obrázek 6.1: Ukázka nahrávky směsi dvou mluvčích



Obrázek 6.2: První mluvčí ze směsi



Obrázek 6.3: Druhý mluvčí ze směsi

Lze si všimnout, že sečtením signálů separovaných mluvčích na obrázku 6.2 a 6.3 dostaneme přesně signál směsi, což lze vyjádřit vztahem

$$x[t] = \sum_{i=1}^C s_i[t] \quad (6.1)$$

, kde  $x[t] \in \mathbb{R}^{1 \times T}$  je diskrétní signál směsi a  $s_i[t] \in \mathbb{R}^{1 \times T}$ , kde  $i = 1, \dots, C$ , je jeden z  $C$  zdrojů[cite studie str3 vlevo].

[[doplnit info o zero means a jendotkove varianci]]

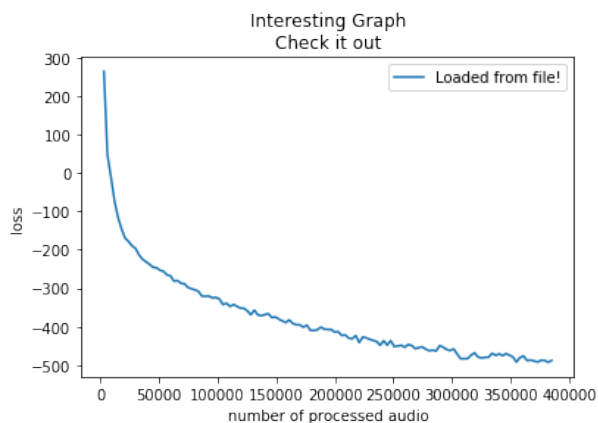
## 6.2 Průběh trénování

[[miniÚvod - konkrétní hyperparams, výkon a čas trénování, seg-len, bylo natrénováno několik modelů ... konkrétní hyperparams, výkon a čas trénování, seg-len]]. Bylo natrénováno několik, různě velkých modelů, které se lišily velikostí jejich hyper-parametrů, konkrétně počtem konvolučních bloků, počtem jejich opakování a velikostí segmentů zpracovávaných nahrávek z datasetu.

[[https://www.tablesgenerator.com/]].

Tabulka 6.1: Hodnoty hyperparametrů trénovaných sítí a jejich přesnost

	X	R	segment length sisnr
2	2	4	8.3
3	2	4	9.45
4	3	4	10.2
8	4	2	13.6



Obrázek 6.4: Příklad grafu loss hodnoty během učení



Obrázek 6.5: Hodnota loss při trénování modelů s různou velikostí hyperparametrů

### 6.3 Experimenty s modely

[[rozdelit tabulku podle HP a spis mensi tabulku ukazujici vysledku pro kon-  
kretni hyper parametry. experimenty: sisnr, reportovat ruzne metriky - stoi,  
pesq, prozkoumat hyper parametry a jak ktere jsou dulezite. Kdyz budu zmen-  
sovat sit tak jak se bude zhorsovat.]]

[[miniÚvod]] Pro vyhodnocení úspěšnosti a kvality modelů byly použity metriky si-snr  
vyjadřující poměr XXX, pesq říkájící YYY a stoi vyjadřující ZZZ [[mozna ref]].

### 6.3.1 Výsledky referenčního modelu TasNet

[[sem nebo do TasNet kapitoly? Jen lehce bych to shrnul ceho dosahli v tech vecich ktere tu merim a vyhodnocuji ja. Pozn: dat to na konec experimentu, muj nej a srovnat s tim jejich nejlepsim, pripadne takovej, kde jsou parametry podobne tem mejm treba.]]

### 6.3.2 Vliv hyper-parametru X

[[zmensovani a zvetsovani X, vysledky testovani, loss pripadne nejake fajnové grafy porovnani vysledky mezi jednotlivymi X, obrázky, grafy]]

### 6.3.3 Vliv hyper-parametru R

[[zmensovani a zvetsovani Y, vysledky testovani, loss pripadne nejake fajnové grafy porovnani vysledky mezi jednotlivymi Y, obrázky, grafy]]

### 6.3.4 Vliv hyper-parametru délka-segmentu

[[zmensovani a zvetsovani SL, vysledky testovani, loss pripadne nejake fajnové grafy porovnani vysledky mezi jednotlivymi SL, obrázky, grafy]]

### 6.3.5 Vyhodnocení a porovnání modelů

[[shrnutí a porovnani jednotlivych zmen mezi sebou komplexne. Pripadne porovnani vysledku s ref studii, obrázky, grafy. Okomentovat, co je tedy dulezity parametr, a ktery jde zmensit a neni moc dulezity z pohledu vysledku evaluace.]]

histogram vysledku skrze jednotlivé promluvy, zjistit pomery jestli je lepsi vysledky muz zena nebo muz muz, atd, v nahrakach je ID, a podle toho to zkosit zmapovat. Zkosit to treba pro nejlepsi model, pripadne pro nejakej jeden mensi.

## 6.4 Možná rozšíření a navrhnutá vylepšení

- variabilnější dataset, mikrofony, šum a bordel prostředí - separace více mluvčích - hlučné prostředí - identifikace konkrétního řečníka - realtime separace poznamka viz backpropagation.

## Kapitola 7

# Závěr

**[[co jak dopadlo, vysledky a vyhodnoceni velikosti modelu a jaky byl nejlepsi]]**

Cílem práce bylo implementovat síť podle architektury TasNet pro separaci mluvčích v časové doméně a porovnat vliv velikosti sítě na kvalitu separace. Síť byla implementována za pomoci frameworku PyTorch a jazyku Python a natrénována na datasetu obsahujícím jednokanálové směsi dvou mluvčích. Trénování proběhlo na **[[X]]** modelech, které se od sebe lišily počtem opakujících se konvolučních bloků, velikostí časové dilatace a délkou vstupních segmentů směsí. Pro účel vyhodnocení modelů byla použita metrika si-snr, která udává poměr chtěného signálu ku šumu na pozadí, tedy obecně kvalitu separace.

Experimenty ukázaly, že během testování nejlépe dopadla síť, která měla 8 konvolučních bloků po 4 opakováních, s délkou vstupního segmentu  $L = 2$  sekundy. Tento model dosáhl po 100 epochách trénování hodnoty až **[[13,4]]** a tím se stal nejúspěšnějším modelem. Při fyzickém poslechu separovaných nahrávek bychom neslyšeli téměř žádný náznak druhého mluvčího. Oproti tomu, nejméně přesný model měl pouze 4 konvoluční bloky, 2 opakování a při délce segmentů  $L = 4$  sekundy dosahoval hodnoty SDR pouze **[[9.2]]**.

Zkoušel jsem separovat také nahrávky, které byly úplně mimo dataset, ale výsledek se nedá hodnotit jako úspěšný, jelikož hraje velkou roli prostředí, mikrofon, šum v pozadí a další vlivy, na které byla neuronová síť naučena. Tento problém by se dal překonat rozšířením trénovacího datasetu o větší škálu nahrávek mluvčích, které by byly pořízeny z různých zařízení v různě rušném prostředí.

**[[Doplnit ještě neco eh]]** Možna.

# Literatura

- [1] BA, J. L., KIROS, J. R. a HINTON, G. E. *Layer Normalization*. 2016.
- [2] BENGIO, Y., SIMARD, P. a FRASCONI, P. Learning long-term dependencies with gradient descent is difficult. *IEEE Transactions on Neural Networks*. 1994, sv. 5, č. 2, s. 157–166.
- [3] CHEN, S. a CHENG, M.-J. Building an Adaptive Evaluation System: A Design Education Application. *Computer-Aided Design & Applications*. Leden 2006, sv. 3, s. 49–58. DOI: 10.1080/16864360.2006.10738441.
- [4] GLOROT, X. a BENGIO, Y. Understanding the difficulty of training deep feedforward neural networks. In: TEH, Y. W. a TITTERINGTON, M., ed. *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*. Chia Laguna Resort, Sardinia, Italy: PMLR, 13–15 May 2010, sv. 9, s. 249–256. Proceedings of Machine Learning Research. Dostupné z: <http://proceedings.mlr.press/v9/glorot10a.html>.
- [5] HE, K., ZHANG, X., REN, S. a SUN, J. *Deep Residual Learning for Image Recognition*. 2015.
- [6] HE, K., ZHANG, X., REN, S. a SUN, J. *Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification*. 2015.
- [7] IAN GOODFELLOW, A. C. *DEEP LEARNING*. MIT Press, 2017. ISBN 9780262035613.
- [8] IOFFE, S. a SZEGEDY, C. *Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift*. 2015.
- [9] KELLEHER, J. D. *DEEP LEARNING / John D. Kelleher*. MIT Press, 2019. ISBN 9780262537551.
- [10] KRIZHEVSKY, A., SUTSKEVER, I. a HINTON, G. E. ImageNet Classification with Deep Convolutional Neural Networks. In: PEREIRA, F., BURGESS, C. J. C., BOTTOU, L. a WEINBERGER, K. Q., ed. *Advances in Neural Information Processing Systems 25*. Curran Associates, Inc., 2012, s. 1097–1105. Dostupné z: <http://papers.nips.cc/paper/4824-imagenet-classification-with-deep-convolutional-neural-networks.pdf>.
- [11] LECUN, Y., BOSER, B., DENKER, J. S., HENDERSON, D., HOWARD, R. E. et al. Backpropagation Applied to Handwritten Zip Code Recognition. *Neural Computation*. 1989, sv. 1, č. 4, s. 541–551.

- [12] LUO, Y. a MESGARANI, N. *Conv-TasNet: Surpassing Ideal Time-Frequency Magnitude Masking for Speech Separation*. 2018.
- [13] NEAL, R. M. Connectionist learning of belief networks. *Artificial Intelligence*. 1992, sv. 56, č. 1, s. 71 – 113. DOI: [https://doi.org/10.1016/0004-3702\(92\)90065-6](https://doi.org/10.1016/0004-3702(92)90065-6). ISSN 0004-3702. Dostupné z: <http://www.sciencedirect.com/science/article/pii/0004370292900656>.
- [14] NWANKPA, C., IJOMAH, W., GACHAGAN, A. a MARSHALL, S. *Activation Functions: Comparison of trends in Practice and Research for Deep Learning*. 2018.
- [15] PEKEL, E. a KARA, S. A COMPREHENSIVE REVIEW FOR ARTIFICIAL NEURAL NETWORK APPLICATION TO PUBLIC TRANSPORTATION. *Sigma Journal of Engineering and Natural Sciences*. Březen 2017, sv. 35, s. 157–179.
- [16] QIAN, W. C. C. X. e. a. Past review, current progress, and challenges ahead on the cocktail party problem. *Frontiers Inf Technol Electronic Eng* 19. Únor 2019, s. 40–63. Dostupné z: <https://doi-org.ezproxy.lib.vutbr.cz/10.1631/FITEE.1700814>.
- [17] RIX, B. J. G. H. M. P. a HEKSTRA, A. P. Perceptual evaluation of speech quality (pesq)-anew method for speech quality assessment of telephonenetworks and codecs. *Acoustics, Speech, and SignalProcessing*. 2001, sv. 2, s. 40–63. Rix, A. W., Beerends, J. G., Hollier, M. P., and Hekstra, A. P. Perceptual evaluation of speech quality (pesq)-anew method for speech quality assessment of telephonenetworks and codecs. In *Acoustics, Speech, and SignalProcessing*, 2001. Proceedings.(ICASSP'01). 2001 IEEE International Conference on, volume 2, pp. 749–752. IEEE, 2001.
- [18] RUDER, S. *An overview of gradient descent optimization algorithms*. 2016.
- [19] TAAL, C. H., HENDRIKS, R. C., HEUSDENS, R. a JENSEN, J. A short-time objective intelligibility measure for time-frequency weighted noisy speech. In: *2010 IEEE International Conference on Acoustics, Speech and Signal Processing*. 2010, s. 4214–4217.
- [20] TAAL, C. H., HENDRIKS, R. C., HEUSDENS, R. a JENSEN, J. An Algorithm for Intelligibility Prediction of Time–Frequency Weighted Noisy Speech. *IEEE Transactions on Audio, Speech, and Language Processing*. 2011, sv. 19, č. 7, s. 2125–2136.
- [21] VENKATARAMANI, S., HIGA, R. a SMARAGDIS, P. *Performance Based Cost Functions for End-to-End Speech Separation*. 2018.
- [22] VINCENT, E., GRIBONVAL, R. a FÉVOTTE, C. Performance measurement in blind audio source separation. *IEEE Transactions on Audio, Speech and Language Processing*. Institute of Electrical and Electronics Engineers. 2006, sv. 14, č. 4, s. 1462–1469. Dostupné z: <https://hal.inria.fr/inria-00544230>.
- [23] WANG, D. a CHEN, J. Supervised Speech Separation Based on Deep Learning: An Overview. *IEEE/ACM Trans. Audio, Speech and Lang. Proc.* IEEE Press. říjen 2018, sv. 26, č. 10, s. 1702–1726. DOI: 10.1109/TASLP.2018.2842159. ISSN 2329-9290. Dostupné z: <https://doi.org/10.1109/TASLP.2018.2842159>.



- [24] ZAEEMZADEH, A., RAHNAVARD, N. a SHAH, M. *Norm-Preservation: Why Residual Networks Can Become Extremely Deep?* 2018.
- [25] ZEILER, M., RANZATO, M., MONGA, R., MAO, M., YANG, K. et al. On Rectified Linear Units For Speech Processing. In: *38th International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. Vancouver: [b.n.], 2013.