

# Úvod do počítačového zpracování řeči

Luděk Bártek

Fakulta infomatiky  
Masarykova univerzita

podzim 2014

# Obsah

- 1 Syntéza řeči – postprocessing
- 2 Značkování prozódie

# Prozodie

- Výstupem syntézy je monotóní hlas bez intonace a přízvuku – zní nepřirozeně
- Doplnění prozodie
  - základní prozodické prvky:
    - výška
    - hlasitost
    - doba trvání
  - nositelem je slabika
  - Větná intonace (prozodie) – závisí na typu věty:
    - otázky zjišťovací (odpověď ano/ne) – rostoucí
    - oznamovací, tázací doplňovací, rozkazovací – klesající
    - řeší se modulací  $F_0$
  - Doplnění přízvuku/důrazu
    - modifikace  $F_0$  a intenzity
    - lokální modifikace větné melodie

## Prozodie – ukázky větné intonace

- Originální promluva (data/masse.wav)
- Oznamovací věta (data/masse-ozn.wav)
- Otázka zjišťovací (data/masse-dotaz.wav)

# Výška základního tónu

- Výška základního tónu odpovídá formantu  $F_0$ .
- Průběh  $F_0$  na vokalickém jádru bývá nelineární.
- Změna intonace není pouhou změnou  $F_0$ 
  - nutno modifikovat i vyšší formanty.
- Na základě důležitosti  $F_0$  se jazyky dělí na:
  - tónové (čínština, vietnamština, ...)
    - čínské slovo -ma- v závislosti na průběhu  $F_0$  může znamenat matka, konopí, kůň, nadávat
  - jazyky s melodickým přízvukem (srbština, slovinština, litevština, norština, švédština, ...)

## Další prozodické vlastnosti

- Intenzita (hlasitost):
  - fyzikální pohled – intenzita signálu v daném časovém okamžiku
  - fyziologický pohled – reakce vnitřního ucha (cortiho ústrojí) na vnímaný zvuk.
  - Tato hlediska se různí.
    - Subjektivní vnímání zvuku neodpovídá ani v prvním přiblížení fyzikální intenzitě signálu.
- Doba trvání:
  - Slabika může mít různou dobu trvání v různém kontextu.
  - Drobné odchylky mohou být i ve stejném kontextu.
  - Typická doba trvání slabiky 50 — 200 milisekund.

## Další prozodické vlastnosti

- Kvalita hlasu
  - chvění hlasu (jitter)
  - nepravidelné výchylky v amplitudě  $F_0$  (shimmer)
  - zbarvení tónu
  - ochraptělost
  - níra znělosti
  - ...
- Rychlost řeči
  - Lze chápat jako převrácenou hodnotu průměrné délky slabiky
  - Lze měřit i jinými způsoby:
    - počtem vyslovených textových znaků za jednotku času (vyhodnocování syntetizérů řeči).

# Další prozodické vlastnosti

## Pokračování

- Pauza
  - tichá
  - vyplněná – obsahuje nějaký charakteristický zvuk (např. eeh)
    - ztížená detekce – hlavní format je blízký formantům samohlásek "a", "e".
- Zaváhání
  - Přímo vypovídá o pragmatice projevu.
  - Důležitý např. pro modifikaci dialogové strategie u dialogových systémů.
  - Typický případ informace obsažené zejména v prozodické vrstvě jazyka.



# Základní odvozené prozodické vlastnosti

- Rytmus (časování):
  - Prozodický prvek odvozený z dob trvání
    - slabik
    - pauz v daném časovém úseku.
- Slovní přízvuk
  - Je odvozen ze všech základních atributů.
  - Je výrazně jazykově závislý:
    - umístění přízvuku ve slově/přízvučné jednotce
    - míra použití prozodických prostředků k jeho vyjádření zejména použití hlasitosti oproti výšce.
- Větný přístup (intonační centrum):
  - zjednodušeně jde o prozodické zvýraznění jádra výpovědi věty

## Základní odvozené prozodické vlastnosti (2.)

- Intonace
  - nejobecněji – časový průběh zvukového spektra hlasu
  - za určující pro melodii se obvykle považuje základní hlasová frekvence – lze zobrazit grafem v závislosti na čase
    - časová závislost základní hlasové frekvence
  - související terminologie:
    - melodie
    - kadence
    - intonační kadence
    - melodém
    - průběh  $F_0$
- Emotivní zbarvení hlasu
  - projevuje se:
    - rychlými změnami hlasitosti a základní frekvence
  - Často přesahují hranici věty.
  - Detekce je důležitá např. pro dialogové systémy – umožňuje zvolit vhodnou dialogovou strategii.

## Základní odvozené prozodické vlastnosti (3.)

- Emfatický přízvuk
  - Vytvářen emotivním zbarvením hlasu.
  - Vyskytuje se např. ve větách pronesených v situacích s výrazným emocionálním kontextem, např.
    - To je tedy opravdu **neslýchané**.
    - Bolí to jak **čert**.
- Kontrastní přízvuk
  - snaha o zdůraznění slova nebo slabiky v kontrastu s jiným slovem nebo slabikou během promluvy nebo dialogu:
    - "řekl jsem do **Šakvic** ne **Rakvic**"
    - "**byte** ne **bit**"

## Základní odvozené prozodické vlastnosti (4.)

- Opakování
  - prozodický atribut silně svázaný s mluvčím.
  - Opakování bývá často variantou výplňkových částí promluvy – mluvčí si ji často ani neuvědomuje (nezaměňovat s koktáním – porucha řeči).
  - Může se jednat o formu zdůraznění – v krajním případě může být považováno za vadu řeči.
- Výplňkové části
  - kromě výplňkové funkce mohou charakterizovat
    - styl mluvčího: „Byl jsi včera na akci, **viď?**”
  - nářečí resp. slang: „**Vole**, ta včerejší spáňka byla hustá, že **vole?**”

## Základní odvozené prozodické vlastnosti (5.)

- Přerušení:
  - častý jev v mluvené řeči na úrovni:
    - vyšších celků (výpověď/promluva, věta, prozodická fráze, ...)
    - uvnitř slov.
  - Mívá návaznost na další prosodické prvky:
    - zaváhání
    - opakování
    - vyplněnou pauzu
    - ...
  - Zvyšuje obtížnost rozpoznávání mluvené řeči – nutno s ním počítat.
- Korekce částí promluvy:
  - Častý jev a to vzhledem k rozdílným částem.
  - Příčiny vzniku:
    - důsledek přeřeknutí,
    - upřesnění předchozí části promluvy,
    - oprava předchozí části promluvy.
  - Často následuje přerušení nebo další prozodické jevy.

# Prozodické segmenty mluvené řeči

- Prozodické segmenty mluvené řeči:
  - Promluva.
  - Prozodická fráze
    - Skupina slov vytvářející jednotný intonační celek.
    - Představuje základní, z prozodického hlediska kompaktní strukturu.
    - Členění do prozodických frází ve velké míře souvisí se syntaktickou strukturou odpovídající věty.
  - Přízvukový takt
    - skupina slabik podřízená jednomu slovnímu přízvuku.
    - V češtině typicky slovo nebo slovo a jednoslabičné slovo.
  - Slabika

# Standards pro syntézu řeči

- Snaha sjednotit jazyky pro popis promluvy pro řečové syntetizéry.
- Definují značkování postihující:
  - prozódii
    - rychlost řeči
    - $F_0$
    - zdůraznění části promluvy
    - pauzu
    - hlasitost
    - ...
  - mluvčího
    - pohlaví
    - věk
    - ...
  - ...
- Používané standardy:
  - SABLE
  - SSML

# SABLE

- Vývoj započat v 2. polovině 90. let
- aplikace XML/SGML
- snaha o zkombinování 3 značkovacích jazyků pro syntézu řeči:
  - SSML – Speech Synthesis Markup Language (W3C, 1999)
  - STML – Spoken Text Markup Language (CSTR Edinburgh University, Lucent Technologies, 1997)
  - JSML – Java Synthesis Markup Language (Sun Microsystems, 2000)
- SABLE



# SABLE

## Základní značky

- SABLE – kořenová značka
- div – slouží k logickému členění dokumentu (odstavec, věta)
- prozodické:
  - EMPH – zdůraznění části promluvy
  - PITCH – výška promluvy
  - VOLUME – úroveň hlasitosti
  - RATE – rychlost
  - BREAK – pauza
- popis hlasu:
  - SPEAKER – popisuje pohlaví a věk mluvčího
- fonetické
  - PRON – výslovnost – fonetický přepis
  - SAYAS – způsob fonetického přepisu (datum, telefon, url, poštovní adresa, ...)
  - LANGUAGE – jazyk promluvy

## SABLE – ukázka

```
<SABLE>
  <DIV TYPE="paragraph">
    <VOLUME LEVEL="quiet">Šepot.</VOLUME>
    <VOLUME LEVEL="medium">
      <RATE SPEED="fast">Rychlá věta.</RATE>
      <PITCH BASE="+50%">Vysoko posazená věta</PITCH>
    </VOLUME>
  </DIV>
</SABLE>
```

# SSML

- Vývoj započat v koncem 90. let
- součást W3C Voice Browser Activity
- Aktuální verze 1.0 (září 2004)

# SSML

## Základní značky

- kořenový element – speak
- strukturní elementy
  - p – odstavec
  - s – věta
- fonetické:
  - say-as – způsob fonetického přepisu (výslovnosti, datum, telefon, url, číslo, ...)
  - phoneme – fonetický přepis dané promluvy
  - sub – substituce (např. přepis zkratk, ...)
- popis hlasu:
  - voice – popis hlasu, kterým se má text přečíst (pohlaví, věk, ...)
- prozodie:
  - emphasis – zdůraznění částí promluvy
  - break – pauza
  - prosody – ovlivňuje prozodické jevy: výšku, průběh základní frekvence, rychlost, item délka trvání promluvy, hlasitost.

## SSML

## Ukázka

```
<?xml version="1.0"?>  
<speak>  
  <voice gender="female">Female voice.</voice>  
  <voice gender="male">Male voice.</voice>  
  <emphais level="soft">Soft emphasis</emphasis>  
  <p>Speech with 5 seconds <break time="5s"/> break.</p>  
  <prosody volume="+6dB">Speech at double volume.</prosody>  
  <prosody volume="-6dB">Speech at half volume.</prosody>  
</speak>
```