

Syntéza a rozpoznávání řeči

Pavel Cenek, Aleš Horák

E-mail: hales@fi.muni.cz
http://nlp.fi.muni.cz/poc_lingv/

Obsah:

- Syntéza řeči
- Rozpoznávání řeči
- Související technologie

Syntéza řeči

- Text to Speech, TTS
- Konverze textu do mluvené podoby
- V ideálním případě by měla syntetizovaná řeč znít tak, jako kdyby daný text přečetl člověk
- Probíhá obvykle ve 4 fázích
 - Normalizace textu
 - Fonetický přepis
 - Prozodický přepis
 - Akustické modelování

První 3 fáze = NLP modul, čtvrtá fáze = DSP modul

Syntéza řeči

- Text to Speech, TTS
- Konverze textu do mluvené podoby
- V ideálním případě by měla syntetizovaná řeč znít tak, jako kdyby daný text přečetl člověk
- Probíhá obvykle ve 4 fázích
 - Normalizace textu
 - Fonetický přepis
 - Prozodický přepis
 - Akustické modelování

První 3 fáze = NLP modul, čtvrtá fáze = DSP modul

Syntéza řeči

- Text to Speech, TTS
- Konverze textu do mluvené podoby
- V ideálním případě by měla syntetizovaná řeč znít tak, jako kdyby daný text přečetl člověk
- Probíhá obvykle ve 4 fázích
 - Normalizace textu
 - Fonetický přepis
 - Prozodický přepis
 - Akustické modelování

První 3 fáze = NLP modul, čtvrtá fáze = DSP modul

Syntéza řeči

- Text to Speech, TTS
- Konverze textu do mluvené podoby
- V ideálním případě by měla syntetizovaná řeč znít tak, jako kdyby daný text přečetl člověk
- Probíhá obvykle ve 4 fázích
 - Normalizace textu
 - Fonetický přepis
 - Prozodický přepis
 - Akustické modelování

První 3 fáze = NLP modul, čtvrtá fáze = DSP modul

Syntéza řeči

- Text to Speech, TTS
- Konverze textu do mluvené podoby
- V ideálním případě by měla syntetizovaná řeč znít tak, jako kdyby daný text přečetl člověk
- Probíhá obvykle ve 4 fázích
 - Normalizace textu
 - Fonetický přepis
 - Prozodický přepis
 - Akustické modelování

První 3 fáze = NLP modul, čtvrtá fáze = DSP modul

Syntéza řeči

- Text to Speech, TTS
- Konverze textu do mluvené podoby
- V ideálním případě by měla syntetizovaná řeč znít tak, jako kdyby daný text přečetl člověk
- Probíhá obvykle ve 4 fázích
 - Normalizace textu
 - Fonetický přepis
 - Prozodický přepis
 - Akustické modelování

První 3 fáze = NLP modul, čtvrtá fáze = DSP modul

Syntéza řeči

- Text to Speech, TTS
- Konverze textu do mluvené podoby
- V ideálním případě by měla syntetizovaná řeč znít tak, jako kdyby daný text přečetl člověk
- Probíhá obvykle ve 4 fázích
 - Normalizace textu
 - Fonetický přepis
 - Prozodický přepis
 - Akustické modelování

První 3 fáze = NLP modul, čtvrtá fáze = DSP modul

Syntéza řeči

- Text to Speech, TTS
- Konverze textu do mluvené podoby
- V ideálním případě by měla syntetizovaná řeč znít tak, jako kdyby daný text přečetl člověk
- Probíhá obvykle ve 4 fázích
 - Normalizace textu
 - Fonetický přepis
 - Prozodický přepis
 - Akustické modelování

První 3 fáze = NLP modul, čtvrtá fáze = DSP modul

Normalizace textu

- Rozčlenění textu na věty
- Rozvinutí zkratk, měrných jednotek, čísel apod.

“130895”

- číslo
- telefonní číslo
- datum
- ...

Normalizace textu

- Rozčlenění textu na věty
- Rozvinutí zkratk, měrných jednotek, čísel apod.

“130895”

- číslo
- telefonní číslo
- datum
- ...

Normalizace textu

- Rozčlenění textu na věty
- Rozvinutí zkratk, měrných jednotek, čísel apod.

“130895”

- číslo
- telefonní číslo
- datum
- ...

Normalizace textu

- Rozčlenění textu na věty
- Rozvinutí zkratk, měrných jednotek, čísel apod.

“130895”

- číslo
- telefonní číslo
- datum
- ...

Normalizace textu

- Rozčlenění textu na věty
- Rozvinutí zkratk, měrných jednotek, čísel apod.

“130895”

- číslo
- telefonní číslo
- datum
- ...

Fonetický přepis

- Převeďte předzpracovaný text do **fonetické podoby** (tj. do tvaru, který popisuje výslovnost daného textu) **SLiDo**
- **Mezinárodní fonetická abeceda (IPA)** – v češtině cca 40 fonémů
- Fonetický přepis **češtiny** musí zohlednit např.
 - Spodoba znělosti (včela/fčela, dub/đup)
 - Krajové zvyky (např. shoda/zhoda nebo schoda)
- Problémy přináší přepis **cizích vlastních jmen** a **cizích slov** obecně (např. faux pas nebo francouzská vlastní jména)
- Dvě základní metody
 - Fonetický přepis založený na pravidlech (např. pro češtinu funguje dobře)
 - Fonetický přepis pomocí výslovnostních lexikonů
- Obě metody lze kombinovat

Fonetický přepis

- Převeďte předzpracovaný text do **fonetické podoby** (tj. do tvaru, který popisuje výslovnost daného textu) SLiDo
- **Mezinárodní fonetická abeceda (IPA)** – v češtině cca 40 fonémů
- Fonetický přepis **češtiny** musí zohlednit např.
 - Spodoba znělosti (včela/fčela, dub/đup)
 - Krajové zvyky (např. shoda/zhoda nebo schoda)
- Problémy přináší přepis **cizích vlastních jmen** a **cizích slov** obecně (např. faux pas nebo francouzská vlastní jména)
- Dvě základní metody
 - Fonetický přepis založený na pravidlech (např. pro češtinu funguje dobře)
 - Fonetický přepis pomocí výslovnostních lexikonů
- Obě metody lze kombinovat

Fonetický přepis

- Převeďte předzpracovaný text do **fonetické podoby** (tj. do tvaru, který popisuje výslovnost daného textu) **SLiDo**
- **Mezinárodní fonetická abeceda (IPA)** – v češtině cca 40 fonémů
- Fonetický přepis **češtiny** musí zohlednit např.
 - Spodoba znělosti (včela/fčela, dub/ďup)
 - Krajské zvyky (např. shoda/zhoda nebo schoda)
- Problémy přináší přepis **cizích vlastních jmen** a **cizích slov** obecně (např. faux pas nebo francouzská vlastní jména)
- Dvě základní metody
 - Fonetický přepis založený na pravidlech (např. pro češtinu funguje dobře)
 - Fonetický přepis pomocí výslovnostních lexikonů
- Obě metody lze kombinovat

Fonetický přepis

- Převeďte předzpracovaný text do **fonetické podoby** (tj. do tvaru, který popisuje výslovnost daného textu) **SLiDo**
- **Mezinárodní fonetická abeceda (IPA)** – v češtině cca 40 fonémů
- Fonetický přepis **češtiny** musí zohlednit např.
 - Spodoba znělosti (**v**čela/**f**čela, du**b**/du**p**)
 - Krajské zvyky (např. **sh**oda/**zh**oda nebo **sch**oda)
- Problémy přináší přepis **cizích vlastních jmen** a **cizích slov** obecně (např. faux pas nebo francouzská vlastní jména)
- Dvě základní metody
 - Fonetický přepis založený na pravidlech (např. pro češtinu funguje dobře)
 - Fonetický přepis pomocí výslovnostních lexikonů
- Obě metody lze kombinovat

Fonetický přepis

- Převeďte předzpracovaný text do **fonetické podoby** (tj. do tvaru, který popisuje výslovnost daného textu) **SLiDo**
- **Mezinárodní fonetická abeceda (IPA)** – v češtině cca 40 fonémů
- Fonetický přepis **češtiny** musí zohlednit např.
 - Spodoba znělosti (včela/fčela, dub/ďup)
 - Krajové zvyky (např. shoda/zhoda nebo schoda)
- Problémy přináší přepis **cizích vlastních jmen** a **cizích slov** obecně (např. faux pas nebo francouzská vlastní jména)
- Dvě základní metody
 - Fonetický přepis založený na pravidlech (např. pro češtinu funguje dobře)
 - Fonetický přepis pomocí výslovnostních lexikonů
- Obě metody lze kombinovat

Fonetický přepis

- Převeďte předzpracovaný text do **fonetické podoby** (tj. do tvaru, který popisuje výslovnost daného textu) **SLiDo**
- **Mezinárodní fonetická abeceda (IPA)** – v češtině cca 40 fonémů
- Fonetický přepis **češtiny** musí zohlednit např.
 - Spodoba znělosti (včela/fčela, dub/đup)
 - Krajové zvyky (např. shoda/zhoda nebo schoda)
- Problémy přináší přepis **cizích vlastních jmen** a **cizích slov** obecně (např. faux pas nebo francouzská vlastní jména)
- Dvě základní metody
 - Fonetický přepis založený na pravidlech (např. pro češtinu funguje dobře)
 - Fonetický přepis pomocí výslovnostních lexikonů
- Obě metody lze kombinovat

Fonetický přepis

- Převeďte předzpracovaný text do **fonetické podoby** (tj. do tvaru, který popisuje výslovnost daného textu) **SLiDo**
- **Mezinárodní fonetická abeceda (IPA)** – v češtině cca 40 fonémů
- Fonetický přepis **češtiny** musí zohlednit např.
 - Spodoba znělosti (včela/fčela, dub/ďup)
 - Krajevé zvyky (např. shoda/zhoda nebo schoda)
- Problémy přináší přepis **cizích vlastních jmen** a **cizích slov** obecně (např. faux pas nebo francouzská vlastní jména)
- Dvě základní metody
 - Fonetický přepis založený **na pravidlech** (např. pro češtinu funguje dobře)
 - Fonetický přepis pomocí **výslovnostních lexikonů**
- Obě metody lze kombinovat

Fonetický přepis

- Převeďte předzpracovaný text do **fonetické podoby** (tj. do tvaru, který popisuje výslovnost daného textu) SLiDo
- **Mezinárodní fonetická abeceda (IPA)** – v češtině cca 40 fonémů
- Fonetický přepis **češtiny** musí zohlednit např.
 - Spodoba znělosti (včela/fčela, dub/đup)
 - Krajové zvyky (např. shoda/zhoda nebo schoda)
- Problémy přináší přepis **cizích vlastních jmen** a **cizích slov** obecně (např. faux pas nebo francouzská vlastní jména)
- Dvě základní metody
 - Fonetický přepis založený **na pravidlech** (např. pro češtinu funguje dobře)
 - Fonetický přepis pomocí **výslovnostních lexikonů**
- Obě metody lze kombinovat

Fonetický přepis

- Převeďte předzpracovaný text do **fonetické podoby** (tj. do tvaru, který popisuje výslovnost daného textu) Slido
- **Mezinárodní fonetická abeceda (IPA)** – v češtině cca 40 fonémů
- Fonetický přepis **češtiny** musí zohlednit např.
 - Spodoba znělosti (včela/fčela, dub/ďup)
 - Krajové zvyky (např. shoda/zhoda nebo schoda)
- Problémy přináší přepis **cizích vlastních jmen** a **cizích slov** obecně (např. faux pas nebo francouzská vlastní jména)
- Dvě základní metody
 - Fonetický přepis založený **na pravidlech** (např. pro češtinu funguje dobře)
 - Fonetický přepis pomocí **výslovnostních lexikonů**
- Obě metody lze kombinovat

Fonetický přepis

- Převeďte předzpracovaný text do **fonetické podoby** (tj. do tvaru, který popisuje výslovnost daného textu) SLiDo
- **Mezinárodní fonetická abeceda (IPA)** – v češtině cca 40 fonémů
- Fonetický přepis **češtiny** musí zohlednit např.
 - Spodoba znělosti (včela/fčela, dub/ďup)
 - Krajové zvyky (např. shoda/zhoda nebo schoda)
- Problémy přináší přepis **cizích vlastních jmen** a **cizích slov** obecně (např. faux pas nebo francouzská vlastní jména)
- Dvě základní metody
 - Fonetický přepis založený **na pravidlech** (např. pro češtinu funguje dobře)
 - Fonetický přepis pomocí **výslovnostních lexikonů**
- Obě metody lze kombinovat

Prozodický přepis

- tzv. **suprasegmentální rysy**
- popisuje řečový proud spolu s přepisem do fonémů
- obohacení textu o informace (viz SSML dále) o **lokálních fyzikálních charakteristikách** výsledné zvukové vlny:
 - **délka fonému** – tempo řeči, pauzy
 - **intonace věty** – vzor pro hladinu základní frekvence (*pitch*)
 - **tón** – v některých (tzv. **tónových**) jazycích určuje význam
 - **lexikální přízvuk** – v **přízvukových** jazycích ovlivňuje délku, hlasitost a tón slov
- kvalitní výpočet prozodie = **přirozenost** syntetizované řeči
např. u *tonálních jazyků* silně ovlivní i porozumění
- Emoce
 - člověk je při projevu používá
 - výzkum syntézi s emocemi je o dost složitější

SSML

Prozodický přepis

- tzv. **suprasegmentální rysy**
- popisuje řečový proud spolu s přepisem do fonémů
- obohacení textu o informace (viz SSML dále) o **lokálních fyzikálních charakteristikách** výsledné zvukové vlny:
 - **délka fonému** – tempo řeči, pauzy
 - **intonace věty** – vzor pro hladinu základní frekvence (*pitch*)
 - **tón** – v některých (tzv. **tónových**) jazycích určuje význam
 - **lexikální přízvuk** – v **přízvukových** jazycích ovlivňuje délku, hlasitost a tón slov
- kvalitní výpočet prozodie = **přirozenost** syntetizované řeči
např. u *tonálních jazyků* silně ovlivní i porozumění
- Emoce
 - člověk je při projevu používá
 - výzkum syntézi s emocemi je o dost složitější

SliDe

Prozodický přepis

- tzv. **suprasegmentální rysy**
- popisuje řečový proud spolu s přepisem do fonémů
- obohacení textu o informace (viz SSML dále) o **lokálních fyzikálních charakteristikách** výsledné zvukové vlny:
 - **délka** fonému – **tempo** řeči, pauzy
 - **intonace** věty – vzor pro hladinu **základní frekvence** (*pitch*)
 - **tón** – v některých (tzv. **tónových**) jazycích určuje význam
 - lexikální **přízvuk** – v **přízvukových jazycích** ovlivňuje délku, hlasitost a tón slov
- kvalitní výpočet prozodie = **přirozenost** syntetizované řeči
např. u *tonálních jazyků* silně ovlivní i porozumění
- Emoce
 - člověk je při projevu používá
 - výzkum syntézi s emocemi je o dost složitější

SLiDo

Prozodický přepis

- tzv. **suprasegmentální rysy**
- popisuje řečový proud spolu s přepisem do fonémů
- obohacení textu o informace (viz SSML dále) o **lokálních fyzikálních charakteristikách** výsledné zvukové vlny:
 - **délka** fonému – **tempo** řeči, pauzy
 - **intonace** věty – vzor pro hladinu **základní frekvence** (*pitch*)
 - **tón** – v některých (tzv. **tónových**) jazycích určuje význam
 - lexikální **přízvuk** – v **přízvukových jazycích** ovlivňuje délku, hlasitost a tón slov
- kvalitní výpočet prozodie = **přirozenost** syntetizované řeči
např. u *tonálních jazyků* silně ovlivní i porozumění
- Emoce
 - člověk je při projevu používá
 - výzkum syntézi s emocemi je o dost složitější

SLiDo

Prozodický přepis

- tzv. **suprasegmentální rysy**
- popisuje řečový proud spolu s přepisem do fonémů
- obohacení textu o informace (viz SSML dále) o **lokálních fyzikálních charakteristikách** výsledné zvukové vlny:
 - **délka** fonému – **tempo** řeči, pauzy
 - **intonace** věty – vzor pro hladinu **základní frekvence** (*pitch*)
 - **tón** – v některých (tzv. **tónových**) jazycích určuje význam
 - lexikální **přízvuk** – v **přízvukových jazycích** ovlivňuje délku, hlasitost a tón slov
- kvalitní výpočet prozodie = **přirozenost** syntetizované řeči
např. u *tonálních jazyků* silně ovlivní i porozumění
- Emoce
 - člověk je při projevu používá
 - výzkum syntézi s emocemi je o dost složitější

SLiDo

Prozodický přepis

- tzv. **suprasegmentální rysy**
- popisuje řečový proud spolu s přepisem do fonémů
- obohacení textu o informace (viz SSML dále) o **lokálních fyzikálních charakteristikách** výsledné zvukové vlny:
 - **délka** fonému – **tempo** řeči, pauzy
 - **intonace** věty – vzor pro hladinu **základní frekvence** (*pitch*)
 - **tón** – v některých (tzv. **tónových**) jazycích určuje význam
 - lexikální **přízvuk** – v **přízvukových jazycích** ovlivňuje délku, hlasitost a tón slov
- kvalitní výpočet prozodie = **přirozenost** syntetizované řeči
např. u *tonálních jazyků* silně ovlivní i porozumění
- Emoce
 - člověk je při projevu používá
 - výzkum syntézi s emocemi je o dost složitější

SliDo

Prozodický přepis

- tzv. **suprasegmentální rysy**
- popisuje řečový proud spolu s přepisem do fonémů
- obohacení textu o informace (viz SSML dále) o **lokálních fyzikálních charakteristikách** výsledné zvukové vlny:
 - **délka** fonému – **tempo** řeči, pauzy
 - **intonace** věty – vzor pro hladinu **základní frekvence** (*pitch*)
 - **tón** – v některých (tzv. **tónových**) jazycích určuje význam
 - lexikální **přízvuk** – v **přízvukových jazycích** ovlivňuje délku, hlasitost a tón slov
- kvalitní výpočet prozodie = **přirozenost** syntetizované řeči
např. u *tonálních jazyků* silně ovlivní i porozumění
- Emoce
 - člověk je při projevu používá
 - výzkum syntézi s emocemi je o dost složitější

SliDo

Prozodický přepis

- tzv. **suprasegmentální rysy**
- popisuje řečový proud spolu s přepisem do fonémů
- obohacení textu o informace (viz SSML dále) o **lokálních fyzikálních charakteristikách** výsledné zvukové vlny:
 - **délka** fonému – **tempo** řeči, pauzy
 - **intonace** věty – vzor pro hladinu **základní frekvence** (*pitch*)
 - **tón** – v některých (tzv. **tónových**) jazycích určuje význam
 - lexikální **přízvuk** – v **přízvukových jazycích** ovlivňuje délku, hlasitost a tón slov
- kvalitní výpočet prozodie = **přirozenost** syntetizované řeči
např. u *tonálních jazyků* silně ovlivní i porozumění
- Emoce
 - člověk je při projevu používá
 - výzkum syntézi s emocemi je o dost složitější

SliDo

Prozodický přepis

- tzv. **suprasegmentální rysy**
- popisuje řečový proud spolu s přepisem do fonémů
- obohacení textu o informace (viz SSML dále) o **lokálních fyzikálních charakteristikách** výsledné zvukové vlny:
 - **délka** fonému – **tempo** řeči, pauzy
 - **intonace** věty – vzor pro hladinu **základní frekvence** (*pitch*)
 - **tón** – v některých (tzv. **tónových**) jazycích určuje význam
 - lexikální **přízvuk** – v **přízvukových jazycích** ovlivňuje délku, hlasitost a tón slov **SLiDo**
- kvalitní výpočet prozodie = **přirozenost** syntetizované řeči
např. u *tonálních jazyků* silně ovlivní i porozumění
- Emoce
 - člověk je při projevu používá
 - výzkum syntézi s emocemi je o dost složitější

Speech Synthesis Markup Language (SSML)

- Doporučení W3C (jako HTML, XML, ...) – standardní způsob pro doplnění fonetiky a prozodie do textu
- Pokrývá první 3 fáze syntézy řeči (normalizace, fonetický přepis, prozodie)
- `<say-as>` – explicitní určení typu dat (např. `Type="Acronym"`, viz Normalizace)
- `<phoneme>` – fonetický přepis textu
- `<voice>` – změna hlasu (atributy *věk*, *muž/žena*, ...)
- `<emphasis>` – přidání/odebrání důrazu
- `<break>` – vložení/zrušení pauzy
- `<prosody>` – ovlivnění prozodie (výška hlasu, kontura, rychlost, hlasitost atd.)

Speech Synthesis Markup Language (SSML)

- Doporučení W3C (jako HTML, XML, ...) – standardní způsob pro doplnění fonetiky a prozodie do textu
- Pokrývá první 3 fáze syntézy řeči (normalizace, fonetický přepis, prozodie)
- `<say-as>` – explicitní určení typu dat (např. `Type="Acronym"`, viz Normalizace)
- `<phoneme>` – fonetický přepis textu
- `<voice>` – změna hlasu (atributy *věk*, *muž/žena*, ...)
- `<emphasis>` – přidání/odebrání důrazu
- `<break>` – vložení/zrušení pauzy
- `<prosody>` – ovlivnění prozodie (výška hlasu, kontura, rychlost, hlasitost atd.)

Speech Synthesis Markup Language (SSML)

- Doporučení W3C (jako HTML, XML, ...) – standardní způsob pro doplnění fonetiky a prozodie do textu
- Pokrývá první 3 fáze syntézy řeči (normalizace, fonetický přepis, prozodie)
- `<say-as>` – explicitní určení typu dat (např. `Type="Acronym"`, viz Normalizace)
- `<phoneme>` – fonetický přepis textu
- `<voice>` – změna hlasu (atributy *věk*, *muž/žena*, ...)
- `<emphasis>` – přidání/odebrání důrazu
- `<break>` – vložení/zrušení pauzy
- `<prosody>` – ovlivnění prozodie (výška hlasu, kontura, rychlost, hlasitost atd.)

Speech Synthesis Markup Language (SSML)

- Doporučení W3C (jako HTML, XML, ...) – standardní způsob pro doplnění fonetiky a prozodie do textu
- Pokrývá první 3 fáze syntézy řeči (normalizace, fonetický přepis, prozodie)
- `<say-as>` – explicitní určení typu dat (např. `Type="Acronym"`, viz Normalizace)
- `<phoneme>` – fonetický přepis textu
- `<voice>` – změna hlasu (atributy *věk*, *muž/žena*, ...)
- `<emphasis>` – přidání/odebrání důrazu
- `<break>` – vložení/zrušení pauzy
- `<prosody>` – ovlivnění prozodie (výška hlasu, kontura, rychlost, hlasitost atd.)

Speech Synthesis Markup Language (SSML)

- Doporučení W3C (jako HTML, XML, ...) – standardní způsob pro doplnění fonetiky a prozodie do textu
- Pokrývá první 3 fáze syntézy řeči (normalizace, fonetický přepis, prozodie)
- `<say-as>` – explicitní určení typu dat (např. `Type="Acronym"`, viz Normalizace)
- `<phoneme>` – fonetický přepis textu
- `<voice>` – změna hlasu (atributy *věk*, *muž/žena*, ...)
- `<emphasis>` – přidání/odebrání důrazu
- `<break>` – vložení/zrušení pauzy
- `<prosody>` – ovlivnění prozodie (výška hlasu, kontura, rychlost, hlasitost atd.)

Speech Synthesis Markup Language (SSML)

- Doporučení W3C (jako HTML, XML, ...) – standardní způsob pro doplnění fonetiky a prozodie do textu
- Pokrývá první 3 fáze syntézy řeči (normalizace, fonetický přepis, prozodie)
- `<say-as>` – explicitní určení typu dat (např. `Type="Acronym"`, viz Normalizace)
- `<phoneme>` – fonetický přepis textu
- `<voice>` – změna hlasu (atributy *věk*, *muž/žena*, ...)
- `<emphasis>` – přidání/odebrání důrazu
- `<break>` – vložení/zrušení pauzy
- `<prosody>` – ovlivnění prozodie (výška hlasu, kontura, rychlost, hlasitost atd.)

Speech Synthesis Markup Language (SSML)

- Doporučení W3C (jako HTML, XML, ...) – standardní způsob pro doplnění fonetiky a prozodie do textu
- Pokrývá první 3 fáze syntézy řeči (normalizace, fonetický přepis, prozodie)
- `<say-as>` – explicitní určení typu dat (např. `Type="Acronym"`, viz Normalizace)
- `<phoneme>` – fonetický přepis textu
- `<voice>` – změna hlasu (atributy *věk*, *muž/žena*, ...)
- `<emphasis>` – přidání/odebrání důrazu
- `<break>` – vložení/zrušení pauzy
- `<prosody>` – ovlivnění prozodie (výška hlasu, kontura, rychlost, hlasitost atd.)

Speech Synthesis Markup Language (SSML)

- Doporučení W3C (jako HTML, XML, ...) – standardní způsob pro doplnění fonetiky a prozodie do textu
- Pokrývá první 3 fáze syntézy řeči (normalizace, fonetický přepis, prozodie)
- `<say-as>` – explicitní určení typu dat (např. `Type="Acronym"`, viz Normalizace)
- `<phoneme>` – fonetický přepis textu
- `<voice>` – změna hlasu (atributy *věk*, *muž/žena*, ...)
- `<emphasis>` – přidání/odebrání důrazu
- `<break>` – vložení/zrušení pauzy
- `<prosody>` – ovlivnění prozodie (výška hlasu, kontura, rychlost, hlasitost atd.)

Speech Synthesis Markup Language (SSML) – příklad

```
<?xml version="1.0" encoding="UTF-8"?>
<vxml version="2.0" xmlns="http://www.w3.org/2001/vxml">
<form>
  <block>
    <prompt>
      <voice gender="male"><emphasis>Hello</emphasis> Jane.</voice>
      <voice gender="female"><emphasis>Hello</emphasis> Mike,
        how <emphasis>are</emphasis> you?</voice>
      <voice gender="male">I am fine. And how are
        <emphasis>you</emphasis> Jane?</voice>
      <voice gender="female">Not bad.</voice>
      <voice gender="male">OK, Goodbye.</voice>
      <voice gender="female"><emphasis>Goodbye</emphasis>
        Mike.</voice>
    </prompt>
  </block>
</form>
</vxml>
```

Akustické modelování

- **Generování** výsledného akustického **signálu** z předzpracovaného textu (řeší DSP modul)
- Dva základní přístupy
 - syntéza řeči v časové oblasti
 - syntéza řeči ve frekvenční oblasti
- v posledních letech i modelování pomocí **hlubokých neuronových sítí (WaveNet)**
 - ze začátku příliš výpočetně náročné pro aplikace v reálném čase
 - aktuálně používané v Google Assistant
 - umožňuje snadněji trénovat nové hlasy

Akustické modelování

- **Generování** výsledného akustického **signálu** z předzpracovaného textu (řeší DSP modul)
- Dva základní přístupy
 - syntéza řeči v **časové oblasti**
 - syntéza řeči ve **frekvenční oblasti**
- v posledních letech i modelování pomocí **hlubokých neuronových sítí (WaveNet)**
 - ze začátku příliš výpočetně náročné pro aplikace v reálném čase
 - aktuálně používané v Google Assistant
 - umožňuje snadněji trénovat nové hlasy

Akustické modelování

- **Generování** výsledného akustického **signálu** z předzpracovaného textu (řeší DSP modul)
- Dva základní přístupy
 - syntéza řeči v **časové oblasti**
 - syntéza řeči ve **frekvenční oblasti**
- v posledních letech i modelování pomocí **hlubokých neuronových sítí (WaveNet)**
 - ze začátku příliš výpočetně náročné pro aplikace v reálném čase
 - aktuálně používané v Google Assistant
 - umožňuje snadněji trénovat nové hlasy

Akustické modelování

- **Generování** výsledného akustického **signálu** z předzpracovaného textu (řeší DSP modul)
- Dva základní přístupy
 - syntéza řeči v **časové oblasti**
 - syntéza řeči ve **frekvenční oblasti**
- v posledních letech i modelování pomocí **hlubokých neuronových sítí (WaveNet)**
 - ze začátku příliš výpočetně náročné pro aplikace v reálném čase
 - aktuálně používané v Google Assistant
 - umožňuje snadněji trénovat nové hlasy

Akustické modelování

- **Generování** výsledného akustického **signálu** z předzpracovaného textu (řeší DSP modul)
- Dva základní přístupy
 - syntéza řeči v **časové oblasti**
 - syntéza řeči ve **frekvenční oblasti**
- v posledních letech i modelování pomocí **hlubokých neuronových sítí (WaveNet)**
 - ze začátku příliš výpočetně náročné pro aplikace v reálném čase
 - aktuálně používané v **Google Assistant**
 - umožňuje snadněji trénovat **nové hlasy**

Akustické modelování

- **Generování** výsledného akustického **signálu** z předzpracovaného textu (řeší DSP modul)
- Dva základní přístupy
 - syntéza řeči v **časové oblasti**
 - syntéza řeči ve **frekvenční oblasti**
- v posledních letech i modelování pomocí **hlubokých neuronových sítí (WaveNet)**
 - ze začátku příliš výpočetně náročné pro aplikace v reálném čase
 - aktuálně používané v **Google Assistant**
 - umožňuje snadněji trénovat **nové hlasy**

Akustické modelování

- **Generování** výsledného akustického **signálu** z předzpracovaného textu (řeší DSP modul)
- Dva základní přístupy
 - syntéza řeči v **časové oblasti**
 - syntéza řeči ve **frekvenční oblasti**
- v posledních letech i modelování pomocí **hlubokých neuronových sítí (WaveNet)**
 - ze začátku příliš výpočetně náročné pro aplikace v reálném čase
 - aktuálně používané v **Google Assistant**
 - umožňuje snadněji trénovat **nové hlasy**

Akustické modelování

- **Generování** výsledného akustického **signálu** z předzpracovaného textu (řeší DSP modul)
- Dva základní přístupy
 - syntéza řeči v **časové oblasti**
 - syntéza řeči ve **frekvenční oblasti**
- v posledních letech i modelování pomocí **hlubokých neuronových sítí (WaveNet)**
 - ze začátku příliš výpočetně náročné pro aplikace v reálném čase
 - aktuálně používané v **Google Assistant**
 - umožňuje snadněji trénovat **nové hlasy**

I. Syntéza řeči v časové oblasti

- = **konkatenativní syntéza**
- Výsledná řeč se skládá z vybraných, dopředu namluvených **segmentů** řeči (difónů, trifónů, slabik apod.)
- Relativně jednoduché na implementaci
- Nutnost vytvoření rozsáhlé **databáze segmentů** (koartikulace, např. 'á' zní jinak v **táta** a **máma**):
 - difóny – **t á t a**
 - trifóny – **t á t a**
 - kombinace – heterogenní segmenty (někdy difóny, trifóny i celá slova)
- Dochází k **deformaci segmentů** jejich spojováním a aplikací prozodických pravidel – “tajemství” komerčních aplikací

I. Syntéza řeči v časové oblasti

- = **konkatenativní syntéza**
- Výsledná řeč se skládá z vybraných, dopředu namluvených **segmentů** řeči (difónů, trifónů, slabik apod.)
- Relativně jednoduché na implementaci
- Nutnost vytvoření rozsáhlé **databáze segmentů** (koartikulace, např. 'á' zní jinak v **táta** a **máma**):
 - difóny – **t á t a**
 - trifóny – **t á t a**
 - kombinace – heterogenní segmenty (někdy difóny, trifóny i celá slova)
- Dochází k **deformaci segmentů** jejich spojováním a aplikací prozodických pravidel – “tajemství” komerčních aplikací

I. Syntéza řeči v časové oblasti

- = **konkatenativní syntéza**
- Výsledná řeč se skládá z vybraných, dopředu namluvených **segmentů** řeči (difónů, trifónů, slabik apod.)
- Relativně jednoduché na implementaci
- Nutnost vytvoření rozsáhlé **databáze segmentů** (koartikulace, např. 'á' zní jinak v **táta** a **máma**):
 - difóny – **t á t a**
 - trifóny – **t á t a**
 - kombinace – heterogenní segmenty (někdy difóny, trifóny i celá slova)
- Dochází k **deformaci segmentů** jejich spojováním a aplikací prozodických pravidel – “tajemství” komerčních aplikací

I. Syntéza řeči v časové oblasti

- = **konkatenativní syntéza**
- Výsledná řeč se skládá z vybraných, dopředu namluvených **segmentů** řeči (difónů, trifónů, slabik apod.)
- Relativně jednoduché na implementaci
- Nutnost vytvoření rozsáhlé **databáze segmentů** (koartikulace, např. 'á' zní jinak v **táta** a **máma**):
 - difóny – t á t a
 - trifóny – t á t a
 - kombinace – heterogenní segmenty (někdy difóny, trifóny i celá slova)
- Dochází k **deformaci segmentů** jejich spojováním a aplikací prozodických pravidel – “tajemství” komerčních aplikací

I. Syntéza řeči v časové oblasti

- = **konkatenativní syntéza**
- Výsledná řeč se skládá z vybraných, dopředu namluvených **segmentů** řeči (difónů, trifónů, slabik apod.)
- Relativně jednoduché na implementaci
- Nutnost vytvoření rozsáhlé **databáze segmentů** (koartikulace, např. 'á' zní jinak v **táta** a **máma**):
 - difóny – t á t a
 - trifóny – t á t a
 - kombinace – heterogenní segmenty (někdy difóny, trifóny i celá slova)
- Dochází k **deformaci segmentů** jejich spojováním a aplikací prozodických pravidel – “tajemství” komerčních aplikací

I. Syntéza řeči v časové oblasti

- = **konkatenativní syntéza**
- Výsledná řeč se skládá z vybraných, dopředu namluvených **segmentů** řeči (difónů, trifónů, slabik apod.)
- Relativně jednoduché na implementaci
- Nutnost vytvoření rozsáhlé **databáze segmentů** (koartikulace, např. 'á' zní jinak v t**á**ta a m**á**ma):
 - difóny – t **á** t a
 - trifóny – t **á** t a
 - kombinace – heterogenní segmenty (někdy difóny, trifóny i celá slova)
- Dochází k **deformaci segmentů** jejich spojováním a aplikací prozodických pravidel – “tajemství” komerčních aplikací

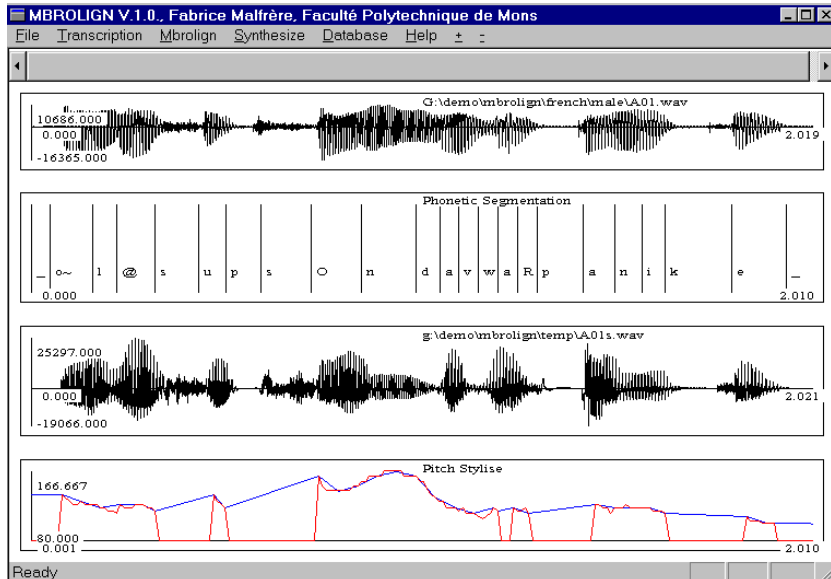
I. Syntéza řeči v časové oblasti

- = **konkatenativní syntéza**
- Výsledná řeč se skládá z vybraných, dopředu namluvených **segmentů** řeči (difónů, trifónů, slabik apod.)
- Relativně jednoduché na implementaci
- Nutnost vytvoření rozsáhlé **databáze segmentů** (koartikulace, např. 'á' zní jinak v t**á**ta a m**á**ma):
 - difóny – t **á** t a
 - trifóny – t **á** t a
 - kombinace – heterogenní segmenty (někdy difóny, trifóny i celá slova)
- Dochází k **deformaci segmentů** jejich spojováním a aplikací prozodických pravidel – “tajemství” komerčních aplikací

I. Syntéza řeči v časové oblasti

- = **konkatenativní syntéza**
- Výsledná řeč se skládá z vybraných, dopředu namluvených **segmentů** řeči (difónů, trifónů, slabik apod.)
- Relativně jednoduché na implementaci
- Nutnost vytvoření rozsáhlé **databáze segmentů** (koartikulace, např. 'á' zní jinak v t**á**ta a m**á**ma):
 - difóny – t **á** t a
 - trifóny – t **á** t a
 - kombinace – heterogenní segmenty (někdy difóny, trifóny i celá slova)
- Dochází k **deformaci segmentů** jejich spojováním a aplikací prozodických pravidel – “tajemství” komerčních aplikací

Semiautomatická tvorba difónové databáze



II. Syntéza řeči ve frekvenční oblasti

2 hlavní přístupy:

- **Modelování hlasového ústrojí**

- generovaný zvuk závisí na parametrech tohoto hlasového ústrojí.
- \oplus velká **flexibilita** (nový hlas lze vytvořit pouhou změnou parametrů)
- \ominus velmi **náročné výpočty** (řeší se fyzikální rovnice modelující situaci ve vokálním traktu, diferenciální rovnice, větš. degradují na válce/koule, ale stejně moc náročné) \Rightarrow v praxi se téměř nepoužívá

- **Formantová syntéza**

modelování (syn) pomocí parametrických popisů hlasového ústrojí
vokálními formantami (f1, f2, f3, f4) - změna hlasového ústrojí se projevuje změnou
polohy těchto formant. (f1 - hluboký hlas, f2 - vysoký hlas, f3 - mužský hlas, f4 - ženský hlas)
závislost na parametrech

II. Syntéza řeči ve frekvenční oblasti

2 hlavní přístupy:

- **Modelování hlasového ústrojí**

- generovaný zvuk závisí na parametrech tohoto hlasového ústrojí.
- ⊕ velká **flexibilita** (nový hlas lze vytvořit pouhou změnou parametrů)
- ⊖ velmi **náročné výpočty** (řeší se fyzikální rovnice modelující situaci ve vokálním traktu, diferenciální rovnice, větš. degradují na válce/koule, ale stejně moc náročné) ⇒ v praxi se téměř nepoužívá

- **Formantová syntéza**

modelování (syn) pomocí formantních parametrů (f0, f1, f2, f3, f4, ...)

→ nejvíce využívaná metoda, protože umožňuje vytvořit hlas velmi podobný lidskému

→ tato metoda využívá fyzikální modely (modely hrdla) a jejich parametry

→ modely hrdla jsou často založeny na fyzikálních principech (např. modely trubek)

→ tyto modely jsou často zjednodušeny, aby byly výpočetně méně náročné

→ modely hrdla jsou často založeny na fyzikálních principech (např. modely trubek)

→ tyto modely jsou často zjednodušeny, aby byly výpočetně méně náročné

II. Syntéza řeči ve frekvenční oblasti

2 hlavní přístupy:

- **Modelování hlasového ústrojí**

- generovaný zvuk závisí na parametrech tohoto hlasového ústrojí.
- \oplus velká **flexibilita** (nový hlas lze vytvořit pouhou změnou parametrů)
- \ominus velmi **náročné výpočty** (řeší se fyzikální rovnice modelující situaci ve vokálním traktu, diferenciální rovnice, větš. degradují na válce/koule, ale stejně moc náročné) \Rightarrow v praxi se téměř nepoužívá

- **Formantová syntéza**

modelování (syn) pomocí formantových filtrů
formantová syntéza (syn) pomocí formantových filtrů
formantová syntéza (syn) pomocí formantových filtrů
formantová syntéza (syn) pomocí formantových filtrů
formantová syntéza (syn) pomocí formantových filtrů
formantová syntéza (syn) pomocí formantových filtrů
formantová syntéza (syn) pomocí formantových filtrů
formantová syntéza (syn) pomocí formantových filtrů
formantová syntéza (syn) pomocí formantových filtrů
formantová syntéza (syn) pomocí formantových filtrů

II. Syntéza řeči ve frekvenční oblasti

2 hlavní přístupy:

- **Modelování hlasového ústrojí**

- generovaný zvuk závisí na parametrech tohoto hlasového ústrojí.
- \oplus velká **flexibilita** (nový hlas lze vytvořit pouhou změnou parametrů)
- \ominus velmi **náročné výpočty** (řeší se fyzikální rovnice modelující situaci ve vokálním traktu, diferenciální rovnice, větš. degradují na válce/koule, ale stejně moc náročné) \Rightarrow v praxi se téměř nepoužívá

- **Formantová syntéza**

II. Syntéza řeči ve frekvenční oblasti

2 hlavní přístupy:

- **Modelování hlasového ústrojí**

- generovaný zvuk závisí na parametrech tohoto hlasového ústrojí.
- \oplus velká **flexibilita** (nový hlas lze vytvořit pouhou změnou parametrů)
- \ominus velmi **náročné výpočty** (řeší se fyzikální rovnice modelující situaci ve vokálním traktu, diferenciální rovnice, větš. degradují na válce/koule, ale stejně moc náročné) \Rightarrow v praxi se téměř nepoužívá

- **Formantová syntéza**

- modelování (jen) hlavních akustických rysů řečového signálu
- model zdroj/filtr – zdroj generuje základní tón pro znělé části řeči a šum pro neznělé části řeči a filtry modifikují zvukové spektrum a napodobují tak hlavní funkce lidského vokálního traktu
- zdroj i filtr jsou řízeny množinou fonetických pravidel \rightarrow syntéza založená na pravidlech
- lze počítat v reálném čase
- mnohem menší data než u konkatenační syntézy \rightarrow vhodné i pro *embedded devices*
- *espeak* – open source projekt espeak.sourceforge.net

II. Syntéza řeči ve frekvenční oblasti

2 hlavní přístupy:

- **Modelování hlasového ústrojí**

- generovaný zvuk závisí na parametrech tohoto hlasového ústrojí.
- \oplus velká **flexibilita** (nový hlas lze vytvořit pouhou změnou parametrů)
- \ominus velmi **náročné výpočty** (řeší se fyzikální rovnice modelující situaci ve vokálním traktu, diferenciální rovnice, větš. degradují na válce/koule, ale stejně moc náročné) \Rightarrow v praxi se téměř nepoužívá

- **Formantová syntéza**

- modelování (jen) **hlavních** akustických rysů řečového signálu
- model **zdroj/filtr** – zdroj generuje základní tón pro znělé části řeči a šum pro neznělé části řeči a filtry modifikují zvukové spektrum a napodobují tak hlavní funkce lidského vokálního traktu
- zdroj i filtr jsou řízeny množinou fonetických **pravidel** \rightarrow syntéza založená na pravidlech
- lze počítat v **reálném čase**
- mnohem **menší data** než u konkatenační syntézy \rightarrow vhodné i pro *embedded devices*
- **espeak** – open source projekt espeak.sourceforge.net

II. Syntéza řeči ve frekvenční oblasti

2 hlavní přístupy:

- **Modelování hlasového ústrojí**

- generovaný zvuk závisí na parametrech tohoto hlasového ústrojí.
- \oplus velká **flexibilita** (nový hlas lze vytvořit pouhou změnou parametrů)
- \ominus velmi **náročné výpočty** (řeší se fyzikální rovnice modelující situaci ve vokálním traktu, diferenciální rovnice, větš. degradují na válce/koule, ale stejně moc náročné) \Rightarrow v praxi se téměř nepoužívá

- **Formantová syntéza**

- modelování (jen) **hlavních** akustických rysů řečového signálu
- model **zdroj/filtr** – zdroj generuje základní tón pro znělé části řeči a šum pro neznělé části řeči a filtry modifikují zvukové spektrum a napodobují tak hlavní funkce lidského vokálního traktu
- zdroj i filtr jsou řízeny množinou fonetických **pravidel** \rightarrow syntéza založená na pravidlech
- lze počítat v **reálném čase**
- mnohem **menší data** než u konkatenační syntézy \rightarrow vhodné i pro *embedded devices*
- **espeak** – open source projekt espeak.sourceforge.net

II. Syntéza řeči ve frekvenční oblasti

2 hlavní přístupy:

- **Modelování hlasového ústrojí**

- generovaný zvuk závisí na parametrech tohoto hlasového ústrojí.
- \oplus velká **flexibilita** (nový hlas lze vytvořit pouhou změnou parametrů)
- \ominus velmi **náročné výpočty** (řeší se fyzikální rovnice modelující situaci ve vokálním traktu, diferenciální rovnice, větš. degradují na válce/koule, ale stejně moc náročné) \Rightarrow v praxi se téměř nepoužívá

- **Formantová syntéza**

- modelování (jen) **hlavních** akustických rysů řečového signálu
- model **zdroj/filtr** – zdroj generuje základní tón pro znělé části řeči a šum pro neznělé části řeči a filtry modifikují zvukové spektrum a napodobují tak hlavní funkce lidského vokálního traktu
- zdroj i filtr jsou řízeny množinou fonetických **pravidel** \rightarrow syntéza založená na pravidlech
- lze počítat v **reálném čase**
- mnohem **menší data** než u konkatentivní syntézy \rightarrow vhodné i pro *embedded devices*
- **espeak** – open source projekt espeak.sourceforge.net

II. Syntéza řeči ve frekvenční oblasti

2 hlavní přístupy:

- **Modelování hlasového ústrojí**

- generovaný zvuk závisí na parametrech tohoto hlasového ústrojí.
- \oplus velká **flexibilita** (nový hlas lze vytvořit pouhou změnou parametrů)
- \ominus velmi **náročné výpočty** (řeší se fyzikální rovnice modelující situaci ve vokálním traktu, diferenciální rovnice, větš. degradují na válce/koule, ale stejně moc náročné) \Rightarrow v praxi se téměř nepoužívá

- **Formantová syntéza**

- modelování (jen) **hlavních** akustických rysů řečového signálu
- model **zdroj/filtr** – zdroj generuje základní tón pro znělé části řeči a šum pro neznělé části řeči a filtry modifikují zvukové spektrum a napodobují tak hlavní funkce lidského vokálního traktu
- zdroj i filtr jsou řízeny množinou fonetických **pravidel** \rightarrow syntéza založená na pravidlech
- lze počítat v **reálném čase**
- mnohem **menší data** než u konkatentivní syntézy \rightarrow vhodné i pro *embedded devices*
- **espeak** – open source projekt espeak.sourceforge.net

II. Syntéza řeči ve frekvenční oblasti

2 hlavní přístupy:

- **Modelování hlasového ústrojí**

- generovaný zvuk závisí na parametrech tohoto hlasového ústrojí.
- \oplus velká **flexibilita** (nový hlas lze vytvořit pouhou změnou parametrů)
- \ominus velmi **náročné výpočty** (řeší se fyzikální rovnice modelující situaci ve vokálním traktu, diferenciální rovnice, větš. degradují na válce/koule, ale stejně moc náročné) \Rightarrow v praxi se téměř nepoužívá

- **Formantová syntéza**

- modelování (jen) **hlavních** akustických rysů řečového signálu
- model **zdroj/filtr** – zdroj generuje základní tón pro znělé části řeči a šum pro neznělé části řeči a filtry modifikují zvukové spektrum a napodobují tak hlavní funkce lidského vokálního traktu
- zdroj i filtr jsou řízeny množinou fonetických **pravidel** \rightarrow syntéza založená na pravidlech
- lze počítat v **reálném čase**
- mnohem **menší data** než u konkatentivní syntézy \rightarrow vhodné i pro *embedded devices*
- **espeak** – open source projekt espeak.sourceforge.net

II. Syntéza řeči ve frekvenční oblasti

2 hlavní přístupy:

- **Modelování hlasového ústrojí**

- generovaný zvuk závisí na parametrech tohoto hlasového ústrojí.
- \oplus velká **flexibilita** (nový hlas lze vytvořit pouhou změnou parametrů)
- \ominus velmi **náročné výpočty** (řeší se fyzikální rovnice modelující situaci ve vokálním traktu, diferenciální rovnice, větš. degradují na válce/koule, ale stejně moc náročné) \Rightarrow v praxi se téměř nepoužívá

- **Formantová syntéza**

- modelování (jen) **hlavních** akustických rysů řečového signálu
- model **zdroj/filtr** – zdroj generuje základní tón pro znělé části řeči a šum pro neznělé části řeči a filtry modifikují zvukové spektrum a napodobují tak hlavní funkce lidského vokálního traktu
- zdroj i filtr jsou řízeny množinou fonetických **pravidel** \rightarrow syntéza založená na pravidlech
- lze počítat v **reálném čase**
- mnohem **menší data** než u konkatentivní syntézy \rightarrow vhodné i pro *embedded devices*
- **espeak** – open source projekt espeak.sourceforge.net

II. Syntéza řeči ve frekvenční oblasti

2 hlavní přístupy:

- **Modelování hlasového ústrojí**

- generovaný zvuk závisí na parametrech tohoto hlasového ústrojí.
- \oplus velká **flexibilita** (nový hlas lze vytvořit pouhou změnou parametrů)
- \ominus velmi **náročné výpočty** (řeší se fyzikální rovnice modelující situaci ve vokálním traktu, diferenciální rovnice, větš. degradují na válce/koule, ale stejně moc náročné) \Rightarrow v praxi se téměř nepoužívá

- **Formantová syntéza**

- modelování (jen) **hlavních** akustických rysů řečového signálu
- model **zdroj/filtr** – zdroj generuje základní tón pro znělé části řeči a šum pro neznělé části řeči a filtry modifikují zvukové spektrum a napodobují tak hlavní funkce lidského vokálního traktu
- zdroj i filtr jsou řízeny množinou fonetických **pravidel** \rightarrow syntéza založená na pravidlech
- lze počítat v **reálném čase**
- mnohem **menší data** než u konkatenační syntézy \rightarrow vhodné i pro *embedded devices*
- **espeak** – open source projekt espeak.sourceforge.net

TTS systémy ve světě

nejčastější použití – telefonní systémy

- ©Nuance (<http://www.nuance.com/>)
- ©Loquendo – nyní Nuance
- ©Acapela group (<http://www.acapela-group.com/>)
 - založena v roce 2004 třemi společnostmi, jedna z nich autor MBROLA
- ©IBM (<http://www.research.ibm.com/tts/>)
- ©AT&T (<http://www.research.att.com/>)
- Festival (<http://www.cstr.ed.ac.uk/projects/festival/>)
- MBROLA (<https://github.com/numediart/MBROLA>)
- FreeTTS (<http://freetts.sourceforge.net/>)

TTS systémy ve světě

nejčastější použití – telefonní systémy

- ©Nuance (<http://www.nuance.com/>)
- ©Loquendo – nyní Nuance
- ©Acapela group (<http://www.acapela-group.com/>)
 - založena v roce 2004 třemi společnostmi, jedna z nich autor MBROLA
- ©IBM (<http://www.research.ibm.com/tts/>)
- ©AT&T (<http://www.research.att.com/>)
- Festival (<http://www.cstr.ed.ac.uk/projects/festival/>)
- MBROLA (<https://github.com/numediart/MBROLA>)
- FreeTTS (<http://freetts.sourceforge.net/>)

TTS systémy ve světě

nejčastější použití – telefonní systémy

- ©Nuance (<http://www.nuance.com/>)
- ©Loquendo – nyní Nuance
- ©Acapela group (<http://www.acapela-group.com/>)
 - založena v roce 2004 třemi společnostmi, jedna z nich autor MBROLA
- ©IBM (<http://www.research.ibm.com/tts/>)
- ©AT&T (<http://www.research.att.com/>)
- Festival (<http://www.cstr.ed.ac.uk/projects/festival/>)
- MBROLA (<https://github.com/numediart/MBROLA>)
- FreeTTS (<http://freetts.sourceforge.net/>)

TTS systémy ve světě

nejčastější použití – telefonní systémy

- ©Nuance (<http://www.nuance.com/>)
- ©Loquendo – nyní Nuance
- ©Acapela group (<http://www.acapela-group.com/>)
 - založena v roce 2004 třemi společnostmi, jedna z nich autor MBROLA
- ©IBM (<http://www.research.ibm.com/tts/>)
- ©AT&T (<http://www.research.att.com/>)
- Festival (<http://www.cstr.ed.ac.uk/projects/festival/>)
- MBROLA (<https://github.com/numediart/MBROLA>)
- FreeTTS (<http://freetts.sourceforge.net/>)

TTS systémy ve světě

nejčastější použití – telefonní systémy

- ©Nuance (<http://www.nuance.com/>)
- ©Loquendo – nyní Nuance
- ©Acapela group (<http://www.acapela-group.com/>)
 - založena v roce 2004 třemi společnostmi, jedna z nich autor MBROLA
- ©IBM (<http://www.research.ibm.com/tts/>)
- ©AT&T (<http://www.research.att.com/>)
- Festival (<http://www.cstr.ed.ac.uk/projects/festival/>)
- MBROLA (<https://github.com/numediart/MBROLA>)
- FreeTTS (<http://freetts.sourceforge.net/>)

TTS systémy ve světě

nejčastější použití – telefonní systémy

- ©Nuance (<http://www.nuance.com/>)
- ©Loquendo – nyní Nuance
- ©Acapela group (<http://www.acapela-group.com/>)
 - založena v roce 2004 třemi společnostmi, jedna z nich autor MBROLA
- ©IBM (<http://www.research.ibm.com/tts/>)
- ©AT&T (<http://www.research.att.com/>)
- Festival (<http://www.cstr.ed.ac.uk/projects/festival/>)
- MBROLA (<https://github.com/numediart/MBROLA>)
- FreeTTS (<http://freetts.sourceforge.net/>)

TTS systémy ve světě

nejčastější použití – telefonní systémy

- ©Nuance (<http://www.nuance.com/>)
- ©Loquendo – nyní Nuance
- ©Acapela group (<http://www.acapela-group.com/>)
 - založena v roce 2004 třemi společnostmi, jedna z nich autor MBROLA
- ©IBM (<http://www.research.ibm.com/tts/>)
- ©AT&T (<http://www.research.att.com/>)
- Festival (<http://www.cstr.ed.ac.uk/projects/festival/>)
- MBROLA (<https://github.com/numediart/MBROLA>)
- FreeTTS (<http://freetts.sourceforge.net/>)

TTS systémy ve světě

nejčastější použití – telefonní systémy

- ©Nuance (<http://www.nuance.com/>)
- ©Loquendo – nyní Nuance
- ©Acapela group (<http://www.acapela-group.com/>)
 - založena v roce 2004 třemi společnostmi, jedna z nich autor MBROLA
- ©IBM (<http://www.research.ibm.com/tts/>)
- ©AT&T (<http://www.research.att.com/>)
- Festival (<http://www.cstr.ed.ac.uk/projects/festival/>)
- MBROLA (<https://github.com/numediart/MBROLA>)
- FreeTTS (<http://freetts.sourceforge.net/>)

České TTS systémy

- EPOS TTS (<http://epos.speech.cz/>) ▶ DEMO
 - Česká akademie věd + Karlova univerzita
- Demosthenes, Popokatepetl
 - LSD FI
- ERIS TTS (<http://www.speechtech.cz/>), heterogenní segmenty ▶ DEMO
 - SpeechTech, s.r.o. + katedra kybernetiky FAV ZČU
 - © verze je nejlepší český
- Český hlas pro MBROLA
 - Mikuláš Piňos, NLP lab FI

České TTS systémy

- EPOS TTS (<http://epos.speech.cz/>) ▶ DEMO
 - Česká akademie věd + Karlova univerzita
- Demosthenes, Popokatepetl
 - LSD FI
- ERIS TTS (<http://www.speechtech.cz/>), heterogenní segmenty ▶ DEMO
 - SpeechTech, s.r.o. + katedra kybernetiky FAV ZČU
 - © verze je nejlepší český
- Český hlas pro MBROLA
 - Mikuláš Piňos, NLP lab FI

České TTS systémy

- EPOS TTS (<http://epos.speech.cz/>) ▶ DEMO
 - Česká akademie věd + Karlova univerzita
- Demosthenes, Popokatepetl
 - LSD FI
- ERIS TTS (<http://www.speechtech.cz/>), heterogenní segmenty ▶ DEMO
 - SpeechTech, s.r.o. + katedra kybernetiky FAV ZČU
© verze je nejlepší český
- Český hlas pro MBROLA
 - Mikuláš Piňos, NLP lab FI

České TTS systémy

- EPOS TTS (<http://epos.speech.cz/>) ▶ DEMO
 - Česká akademie věd + Karlova univerzita
- Demosthenes, Popokatepetl
 - LSD FI
- ERIS TTS (<http://www.speechtech.cz/>), heterogenní segmenty ▶ DEMO
 - SpeechTech, s.r.o. + katedra kybernetiky FAV ZČU
© verze je nejlepší český
- Český hlas pro MBROLA
 - Mikuláš Piňos, NLP lab FI

Obsah

- 1 Syntéza řeči
 - Normalizace textu
 - Fonetický přepis
 - Prozodický přepis
 - SSML
 - Akustické modelování
 - TTS systémy
- 2 Rozpoznávání řeči
 - Porovnávání vektorů parametrů
 - ASR systémy
- 3 Související technologie

Rozpoznávání řeči

- Automatic Speech Recognition, ASR
- Konverze řeči na text
 - Výstupem je většinou množina hypotéz spolu s pravděpodobností správnosti dané hypotézy. K výběru správné hypotézy se běžně využívají jazykové modely
- Lze zhruba rozdělit na
 - Rozpoznávání izolovaných slov – slyšitelná pauza mezi slovy
 - Rozpoznávání kontinuální řeči – plynulá řeč (řeč školeného mluvčího nebo čtený text)
 - Rozpoznávání spontánní řeči – přeroky, pauzy, začátky vět (*false-starts*)

Rozpoznávání řeči

- Automatic Speech Recognition, ASR
- Konverze řeči na text
 - Výstupem je většinou množina hypotéz spolu s pravděpodobností správnosti dané hypotézy. K výběru správné hypotézy se běžně využívají jazykové modely
- Lze zhruba rozdělit na
 - Rozpoznávání izolovaných slov – slyšitelná pauza mezi slovy
 - Rozpoznávání kontinuální řeči – plynulá řeč (řeč školeného mluvčího nebo čtený text)
 - Rozpoznávání spontánní řeči – přeroky, pauzy, začátky vět (*false-starts*)

Rozpoznávání řeči

- Automatic Speech Recognition, ASR
- Konverze řeči na text
 - Výstupem je většinou množina hypotéz spolu s pravděpodobností správnosti dané hypotézy. K výběru správné hypotézy se běžně využívají jazykové modely
- Lze zhruba rozdělit na
 - Rozpoznávání izolovaných slov – slyšitelná pauza mezi slovy
 - Rozpoznávání kontinuální řeči – plynulá řeč (řeč školeného mluvčího nebo čtený text)
 - Rozpoznávání spontánní řeči – přeroky, pauzy, začátky vět (*false-starts*)

Rozpoznávání řeči

- Automatic Speech Recognition, ASR
- Konverze řeči na text
 - Výstupem je většinou množina hypotéz spolu s pravděpodobností správnosti dané hypotézy. K výběru správné hypotézy se běžně využívají jazykové modely
- Lze zhruba rozdělit na
 - Rozpoznávání izolovaných slov – slyšitelná pauza mezi slovy
 - Rozpoznávání kontinuální řeči – plynulá řeč (řeč školeného mluvčího nebo čtený text)
 - Rozpoznávání spontánní řeči – přeroky, pauzy, začátky vět (*false-starts*)

Rozpoznávání řeči

- Automatic Speech Recognition, ASR
- Konverze řeči na text
 - Výstupem je většinou množina hypotéz spolu s pravděpodobností správnosti dané hypotézy. K výběru správné hypotézy se běžně využívají jazykové modely
- Lze zhruba rozdělit na
 - Rozpoznávání izolovaných slov – slyšitelná pauza mezi slovy
 - Rozpoznávání kontinuální řeči – plynulá řeč (řeč školeného mluvčího nebo čtený text)
 - Rozpoznávání spontánní řeči – přeroky, pauzy, začátky vět (*false-starts*)

Rozpoznávání řeči

- Automatic Speech Recognition, ASR
- Konverze řeči na text
 - Výstupem je většinou množina hypotéz spolu s pravděpodobností správnosti dané hypotézy. K výběru správné hypotézy se běžně využívají jazykové modely
- Lze zhruba rozdělit na
 - Rozpoznávání izolovaných slov – slyšitelná pauza mezi slovy
 - Rozpoznávání kontinuální řeči – plynulá řeč (řeč školeného mluvčího nebo čtený text)
 - Rozpoznávání spontánní řeči – přeroky, pauzy, začátky vět (*false-starts*)

Rozpoznávání řeči

- Automatic Speech Recognition, ASR
- Konverze řeči na text
 - Výstupem je většinou množina hypotéz spolu s pravděpodobností správnosti dané hypotézy. K výběru správné hypotézy se běžně využívají jazykové modely
- Lze zhruba rozdělit na
 - Rozpoznávání izolovaných slov – slyšitelná pauza mezi slovy
 - Rozpoznávání kontinuální řeči – plynulá řeč (řeč školeného mluvčího nebo čtený text)
 - Rozpoznávání spontánní řeči – přeroky, pauzy, začátky vět (*false-starts*)

Rozpoznávání řeči pokrač.

- Diktovací stroje (např. Dragon Naturally Speaking)

- Schopné rozpoznat cokoliv
- N -gramové statistické jazykové modely
- Závislé na mluvčím (je potřeba je natrénovat)

SliDo

- Rozpoznávače založené na gramatikách

- Rozpoznají jen fráze popsané (regulární) gramatikou (gramatika = jazykový model)

```
S → "Jedu do " MESTO
MESTO → "Prahy" | "Brna"
```

- Nezávislé na mluvčím – telefonní aplikace
- Speech Recognition Grammar Specification (SRGS)

Standard 192 komparativ, 3 la 2007
 standard 2. vydání – 2007, se 2007
 standard 2. vydání – 2007, se 2007

Rozpoznávání řeči pokrač.

- Diktovací stroje (např. Dragon Naturally Speaking)

- Schopné rozpoznat cokoliv
- N -gramové statistické jazykové modely
- Závislé na mluvcím (je potřeba je natrénovat)

SliDo

- Rozpoznávače založené na gramatikách

- Rozpoznají jen fráze popsané (regulární) gramatikou (gramatika = jazykový model)

$$S \rightarrow "Jedu do " MESTO$$

$$MESTO \rightarrow "Prahy" | "Brna"$$

- Nezávislé na mluvcím – telefonní aplikace
- Speech Recognition Grammar Specification (SRGS)

SRGS: <http://www.siprec.org/specifications/srgs/>

SRGS: <http://www.siprec.org/specifications/srgs/2.0/>

SRGS: <http://www.siprec.org/specifications/srgs/2.0.1/>

Rozpoznávání řeči pokrač.

- Diktovací stroje (např. Dragon Naturally Speaking)

- Schopné rozpoznat cokoliv
- N -gramové statistické jazykové modely
- Závislé na mluvčím (je potřeba je natrénovat)

SLiDo

- Rozpoznávače založené na gramatikách

- Rozpoznají jen fráze popsané (regulární) gramatikou (gramatika = jazykový model)

$$S \rightarrow "Jedu do " MESTO$$

$$MESTO \rightarrow "Prahy" | "Brna"$$

- Nezávislé na mluvčím – telefonní aplikace
- Speech Recognition Grammar Specification (SRGS)

Standard 1920-2000, 3. vydání

Standard 2006-2009, 1. vydání

Standard 2009-2012, 1. vydání

Standard 2012-2015, 1. vydání

Standard 2015-2018, 1. vydání

Standard 2018-2021, 1. vydání

Rozpoznávání řeči pokrač.

- Diktovací stroje (např. Dragon Naturally Speaking)

- Schopné rozpoznat cokoliv
- N -gramové statistické jazykové modely
- Závislé na mluvčím (je potřeba je natrénovat)

SliDo

- Rozpoznávače založené na gramatikách

- Rozpoznají jen fráze popsané (regulární) gramatikou (gramatika = jazykový model)

```
S → "Jedu do " MESTO
MESTO → "Prahy" | "Brna"
```

- Nezávislé na mluvčím – telefonní aplikace
- Speech Recognition Grammar Specification (SRGS)

Standard 1920-2000, 3. vydání

http://www.siprec.org/standards/3rd-edition/srgs.html

Rozpoznávání řeči pokrač.

- Diktovací stroje (např. Dragon Naturally Speaking)
 - Schopné rozpoznat cokoliv
 - N -gramové statistické jazykové modely
 - Závislé na mluvčím (je potřeba je natrénovat)
- Rozpoznávače založené na gramatikách
 - Rozpoznají jen fráze popsané (regulární) gramatikou (gramatika = jazykový model)

$$S \rightarrow \text{"Jedu do " MESTO}$$

$$\text{MESTO} \rightarrow \text{"Prahy" | "Brna"}$$
 - Nezávislé na mluvčím – telefonní aplikace
 - Speech Recognition Grammar Specification (SRGS)
 - standard W3 konzorcia, à la BNF
 - existují 2 notace – XML a šipková pro čtení
 - dá se do ní dát i "význam" vstupu



Rozpoznávání řeči pokrač.

- Diktovací stroje (např. Dragon Naturally Speaking)
 - Schopné rozpoznat cokoliv
 - N -gramové statistické jazykové modely
 - Závislé na mluvčím (je potřeba je natrénovat)
- Rozpoznávače založené na gramatikách
 - Rozpoznají jen fráze popsané (regulární) gramatikou (gramatika = jazykový model)

Slido

$S \rightarrow \text{"Jedu do "MESTO}$
 $\text{MESTO} \rightarrow \text{"Prahy"} \mid \text{"Brna"}$

- Nezávislé na mluvčím – telefonní aplikace
- Speech Recognition Grammar Specification (SRGS)
 - standard W3 konzorcia, à la BNF
 - existují 2 notace – XML a šipková pro čtení
 - dá se do ní dát i "význam" vstupu

Rozpoznávání řeči pokrač.

- Diktovací stroje (např. Dragon Naturally Speaking)
 - Schopné rozpoznat cokoliv
 - N -gramové statistické jazykové modely
 - Závislé na mluvčím (je potřeba je natrénovat)
- Rozpoznávače založené na gramatikách
 - Rozpoznají jen fráze popsané (regulární) gramatikou (gramatika = jazykový model)

$S \rightarrow \text{"Jedu do "MESTO}$
 $\text{MESTO} \rightarrow \text{"Prahy"} \mid \text{"Brna"}$

- Nezávislé na mluvčím – telefonní aplikace
- Speech Recognition Grammar Specification (SRGS)
 - standard W3 konzorcia, à la BNF
 - existují 2 notace – XML a šipková pro čtení
 - dá se do ní dát i "význam" vstupu

Slido

Rozpoznávání řeči pokrač.

- Diktovací stroje (např. Dragon Naturally Speaking)
 - Schopné rozpoznat cokoliv
 - N -gramové statistické jazykové modely
 - Závislé na mluvčím (je potřeba je natrénovat)
- Rozpoznávače založené na gramatikách
 - Rozpoznají jen fráze popsané (regulární) gramatikou (gramatika = jazykový model)

$S \rightarrow \text{"Jedu do "MESTO}$
 $\text{MESTO} \rightarrow \text{"Prahy"} \mid \text{"Brna"}$
 - Nezávislé na mluvčím – telefonní aplikace
 - Speech Recognition Grammar Specification (SRGS)
 - standard W3 konzorcia, à la BNF
 - existují 2 notace – XML a šipková pro čtení
 - dá se do ní dát i “význam” vstupu

Slido

Rozpoznávání řeči pokrač.

- Diktovací stroje (např. Dragon Naturally Speaking)
 - Schopné rozpoznat cokoliv
 - N -gramové statistické jazykové modely
 - Závislé na mluvčím (je potřeba je natrénovat)
- Rozpoznávače založené na gramatikách
 - Rozpoznají jen fráze popsané (regulární) gramatikou (gramatika = jazykový model)

$S \rightarrow \text{"Jedu do "MESTO}$
 $\text{MESTO} \rightarrow \text{"Prahy"} \mid \text{"Brna"}$
 - Nezávislé na mluvčím – telefonní aplikace
 - Speech Recognition Grammar Specification (SRGS)
 - standard W3 konzorcia, à la BNF
 - existují 2 notace – XML a šipková pro čtení
 - dá se do ní dát i “význam” vstupu

Slido

Rozpoznávání řeči pokrač.

- Diktovací stroje (např. Dragon Naturally Speaking)
 - Schopné rozpoznat cokoliv
 - N -gramové statistické jazykové modely
 - Závislé na mluvčím (je potřeba je natrénovat)
- Rozpoznávače založené na gramatikách
 - Rozpoznají jen fráze popsané (regulární) gramatikou (gramatika = jazykový model)

$S \rightarrow \text{"Jedu do "MESTO}$
 $\text{MESTO} \rightarrow \text{"Prahy"} \mid \text{"Brna"}$
 - Nezávislé na mluvčím – telefonní aplikace
 - Speech Recognition Grammar Specification (SRGS)
 - standard W3 konzorcia, à la BNF
 - existují 2 notace – XML a šipková pro čtení
 - dá se do ní dát i "význam" vstupu

Slide

Rozpoznávání řeči pokrač.

- Diktovací stroje (např. Dragon Naturally Speaking)
 - Schopné rozpoznat cokoliv
 - N -gramové statistické jazykové modely
 - Závislé na mluvčím (je potřeba je natrénovat)
- Rozpoznávače založené na gramatikách
 - Rozpoznají jen fráze popsané (regulární) gramatikou (gramatika = jazykový model)

$$S \rightarrow \text{"Jedu do "MESTO}$$
$$\text{MESTO} \rightarrow \text{"Prahy"} \mid \text{"Brna"}$$
 - Nezávislé na mluvčím – telefonní aplikace
 - Speech Recognition Grammar Specification (SRGS)
 - standard W3 konzorcia, à la BNF
 - existují 2 notace – XML a šipková pro čtení
 - dá se do ní dát i "význam" vstupu

Slide

Rozpoznávání řeči pokrač.

Probíhá obvykle ve 3 fázích:

1. Vstup signálu

- Amplituda akustického vlnění je snímána v pravidelných intervalech a uložena ve formě celého čísla (digitalizace a vzorkování signálu)

2. Vytvoření akustických charakteristik signálu (akustické vektory)

- Snižuje variabilitu a odstraňuje redundanci (řeč 300 000 × redundantní)
- Počítají se rozdělením na segmenty 10–40 ms, ze kterých se odečítají charakteristiky, jako je počet průchodů nulou nebo prvních 12 koeficientů FFT (cca 40 čísel, není přesně dané které, ale výběr velice ovlivní výsledek)

3. Porovnávání vektorů parametrů

- K získané sekvenci vektorů parametrů se hledá co nejpodobnější sekvence známých, předem naučených, vektorů reprezentující např. fonémy, trifóny, slabiky, celá slova apod.

Rozpoznávání řeči pokrač.

Probíhá obvykle ve 3 fázích:

1. Vstup signálu

- Amplituda akustického vlnění je snímána v pravidelných intervalech a uložena ve formě celého čísla (digitalizace a vzorkování signálu)

2. Vytvoření akustických charakteristik signálu (akustické vektory)

- Snižuje variabilitu a odstraňuje redundanci (řeč 300 000× redundantní)
- Počítají se rozdělením na segmenty 10–40 ms, ze kterých se odečítají charakteristiky, jako je počet průchodů nulou nebo prvních 12 koeficientů FFT (cca 40 čísel, není přesně dané které, ale výběr velice ovlivní výsledek)

3. Porovnávání vektorů parametrů

- K získané sekvenci vektorů parametrů se hledá co nejpodobnější sekvence známých, předem naučených, vektorů reprezentující např. fonémy, trifóny, slabiky, celá slova apod.

Rozpoznávání řeči pokrač.

Probíhá obvykle ve 3 fázích:

1. Vstup signálu

- Amplituda akustického vlnění je snímána v pravidelných intervalech a uložena ve formě celého čísla (digitalizace a vzorkování signálu)

2. Vytvoření akustických charakteristik signálu (akustické vektory)

- Snižuje variabilitu a odstraňuje redundanci (řeč 300 000× redundantní)
- Počítají se rozdělením na segmenty 10–40 ms, ze kterých se odečítají charakteristiky, jako je počet průchodů nulou nebo prvních 12 koeficientů FFT (cca 40 čísel, není přesně dané které, ale výběr velice ovlivní výsledek)

3. Porovnávání vektorů parametrů

- K získané sekvenci vektorů parametrů se hledá co nejpodobnější sekvence známých, předem naučených, vektorů reprezentující například fonémy, trifóny, slabiky, celá slova apod.

Rozpoznávání řeči pokrač.

Probíhá obvykle ve 3 fázích:

1. Vstup signálu

- Amplituda akustického vlnění je snímána v pravidelných intervalech a uložena ve formě celého čísla (digitalizace a vzorkování signálu)

2. Vytvoření akustických charakteristik signálu (akustické vektory)

- Snižuje variabilitu a odstraňuje redundanci (řeč 300 000× redundantní)
- Počítají se rozdělením na segmenty 10–40 ms, ze kterých se odečítají charakteristiky, jako je počet průchodů nulou nebo prvních 12 koeficientů FFT (cca 40 čísel, není přesně dané které, ale výběr velice ovlivní výsledek)

3. Porovnávání vektorů parametrů

- K získané sekvenci vektorů parametrů se hledá co nejpodobnější sekvence známých, předem naučených, vektorů reprezentující např. fonémy, trifóny, slabiky, celá slova apod.

Rozpoznávání řeči pokrač.

Probíhá obvykle ve 3 fázích:

1. Vstup signálu

- Amplituda akustického vlnění je snímána v pravidelných intervalech a uložena ve formě celého čísla (digitalizace a vzorkování signálu)

2. Vytvoření akustických charakteristik signálu (akustické vektory)

- Snižuje variabilitu a odstraňuje redundanci (řeč 300 000× redundantní)
- Počítají se rozdělením na segmenty 10–40 ms, ze kterých se odečítají charakteristiky, jako je počet průchodů nulou nebo prvních 12 koeficientů FFT (cca 40 čísel, není přesně dané které, ale výběr velice ovlivní výsledek)

3. Porovnávání vektorů parametrů

- K získané sekvenci vektorů parametrů se hledá co nejpodobnější sekvence známých, předem naučených, vektorů reprezentující např. fonémy, trifóny, slabiky, celá slova apod.

Porovnávání vektorů parametrů

- Algoritmus borcení časové osy (dynamic time warping, DTW)
 - odstraňuje časové nerovnoměrnosti v akustickém signálu
- Skryté Markovovy modely (*Hidden Markov Models, HMM*)
 - Pravděpodobnostní konečné automaty
 - V každém okamžiku je hlasové ústrojí v určitém stavu a může s určitou pravděpodobností přejít do jednoho z následujících stavů
 - Jako doplněk se mohou využít neuronové sítě
 - Je nejprve potřeba natrénovat za pomoci dat z řečového korpusu

Porovnávání vektorů parametrů

- Algoritmus borcení časové osy (dynamic time warping, DTW)
 - odstraňuje časové nerovnoměrnosti v akustickém signálu
- Skryté Markovovy modely (*Hidden Markov Models, HMM*)
 - Pravděpodobnostní konečné automaty
 - V každém okamžiku je hlasové ústrojí v určitém stavu a může s určitou pravděpodobností přejít do jednoho z následujících stavů
 - Jako doplněk se mohou využít neuronové sítě
 - Je nejprve potřeba natrénovat za pomoci dat z řečového korpusu

Porovnávání vektorů parametrů

- Algoritmus borcení časové osy (dynamic time warping, DTW)
 - odstraňuje časové nerovnoměrnosti v akustickém signálu
- Skryté Markovovy modely (*Hidden Markov Models, HMM*)
 - Pravděpodobnostní konečné automaty
 - V každém okamžiku je hlasové ústrojí v určitém stavu a může s určitou pravděpodobností přejít do jednoho z následujících stavů
 - Jako doplněk se mohou využít neuronové sítě
 - Je nejprve potřeba natrénovat za pomoci dat z řečového korpusu

Porovnávání vektorů parametrů

- Algoritmus borcení časové osy (dynamic time warping, DTW)
 - odstraňuje časové nerovnoměrnosti v akustickém signálu
- Skryté Markovovy modely (*Hidden Markov Models, HMM*)
 - Pravděpodobnostní konečné automaty
 - V každém okamžiku je hlasové ústrojí v určitém stavu a může s určitou pravděpodobností přejít do jednoho z následujících stavů
 - Jako doplněk se mohou využít neuronové sítě
 - Je nejprve potřeba natrénovat za pomoci dat z řečového korpusu

Porovnávání vektorů parametrů

- Algoritmus borcení časové osy (dynamic time warping, DTW)
 - odstraňuje časové nerovnoměrnosti v akustickém signálu
- Skryté Markovovy modely (*Hidden Markov Models, HMM*)
 - Pravděpodobnostní konečné automaty
 - V každém okamžiku je hlasové ústrojí v určitém stavu a může s určitou pravděpodobností přejít do jednoho z následujících stavů
 - Jako doplněk se mohou využít neuronové sítě
 - Je nejprve potřeba natrénovat za pomoci dat z řečového korpusu

Porovnávání vektorů parametrů

- Algoritmus borcení časové osy (dynamic time warping, DTW)
 - odstraňuje časové nerovnoměrnosti v akustickém signálu
- Skryté Markovovy modely (*Hidden Markov Models, HMM*)
 - Pravděpodobnostní konečné automaty
 - V každém okamžiku je hlasové ústrojí v určitém stavu a může s určitou pravděpodobností přejít do jednoho z následujících stavů
 - Jako doplněk se mohou využít neuronové sítě
 - Je nejprve potřeba natrénovat za pomoci dat z řečového korpusu

ASR systémy ve světě

- ©Nuance (<http://www.nuance.com/>)
- ©Loquendo – nyní Nuance
- ©LumenVox (<http://www.lumenvox.com/>)
- ©IBM ViaVoice – nyní Nuance Dragon Naturally Speaking
<http://www.nuance.com/dragon/>
- Sphinx (<http://cmusphinx.sourceforge.net/>)

ASR systémy ve světě

- ©Nuance (<http://www.nuance.com/>)
- ©Loquendo – nyní Nuance
- ©LumenVox (<http://www.lumenvox.com/>)
- ©IBM ViaVoice – nyní Nuance Dragon Naturally Speaking
<http://www.nuance.com/dragon/>
- Sphinx (<http://cmusphinx.sourceforge.net/>)

ASR systémy ve světě

- ©Nuance (<http://www.nuance.com/>)
- ©Loquendo – nyní Nuance
- ©LumenVox (<http://www.lumenvox.com/>)
- ©IBM ViaVoice – nyní Nuance Dragon Naturally Speaking
<http://www.nuance.com/dragon/>
- Sphinx (<http://cmusphinx.sourceforge.net/>)

ASR systémy ve světě

- ©Nuance (<http://www.nuance.com/>)
- ©Loquendo – nyní Nuance
- ©LumenVox (<http://www.lumenvox.com/>)
- ©IBM ViaVoice – nyní Nuance Dragon Naturally Speaking
<http://www.nuance.com/dragon/>
- Sphinx (<http://cmusphinx.sourceforge.net/>)

ASR systémy ve světě

- ©Nuance (<http://www.nuance.com/>)
- ©Loquendo – nyní Nuance
- ©LumenVox (<http://www.lumenvox.com/>)
- ©IBM ViaVoice – nyní Nuance Dragon Naturally Speaking
<http://www.nuance.com/dragon/>
- Sphinx (<http://cmusphinx.sourceforge.net/>)

České ASR systémy

- Laboratoř počítačového zpracování řeči na Fakultě mechatroniky Technické univerzity v Liberci (<http://www.ite.tul.cz/speechlab/>)
- ERIS ASR (<http://www.speechtech.cz/>)
 - SpeechTech, s.r.o. + katedra kybernetiky FAV ZČU
- Speech@FIT VUT Brno (<http://speech.fit.vutbr.cz/>)
 - keyword spotting – jestli se vyskytlo dané slovo v běžné řeči

České ASR systémy

- Laboratoř počítačového zpracování řeči na Fakultě mechatroniky Technické univerzity v Liberci (<http://www.ite.tul.cz/speechlab/>)
- ERIS ASR (<http://www.speechtech.cz/>)
 - SpeechTech, s.r.o. + katedra kybernetiky FAV ZČU
- Speech@FIT VUT Brno (<http://speech.fit.vutbr.cz/>)
 - keyword spotting – jestli se vyskytlo dané slovo v běžné řeči

České ASR systémy

- Laboratoř počítačového zpracování řeči na Fakultě mechatroniky Technické univerzity v Liberci (<http://www.ite.tul.cz/speechlab/>)
- ERIS ASR (<http://www.speechtech.cz/>)
 - SpeechTech, s.r.o. + katedra kybernetiky FAV ZČU
- Speech@FIT VUT Brno (<http://speech.fit.vutbr.cz/>)
 - keyword spotting – jestli se vyskytlo dané slovo v běžné řeči

Mozilla Common Voice

- voice.mozilla.org
- iniciativa Mozilly pro vytvoření komunitního ASR
- můžete sami přispět k dostupnosti rozpoznávání řeči pro váš jazyk
- uživatelé:
 - vybírají věty (je potřeba > 5,000 vět/jazyk)
 - kontrolují věty
 - nahrávají věty (za všechny jazyky je nyní nahráno 18,000 hodin, za češtinu 67 hodin)
 - kontrolují nahrávky vět
- <https://commonvoice.mozilla.org/cs>,
<https://www.mozilla.cz/zpravicky/tag/voice/>

Mozilla Common Voice

- voice.mozilla.org
- iniciativa Mozilly pro vytvoření **komunitního ASR**
- můžete sami přispět k dostupnosti **rozpoznávání řeči** pro váš jazyk
- uživatelé:
 - vybírají věty (je potřeba > 5,000 vět/jazyk)
 - kontrolují věty
 - nahrávají věty (za všechny jazyky je nyní **nahráno 18,000 hodin**, za češtinu **67 hodin**)
 - kontrolují nahrávky vět
- <https://commonvoice.mozilla.org/cs>,
<https://www.mozilla.cz/zpravicky/tag/voice/>

Mozilla Common Voice

- voice.mozilla.org
- iniciativa Mozilly pro vytvoření **komunitního ASR**
- můžete sami přispět k dostupnosti **rozpoznávání řeči** pro váš jazyk
- uživatelé:
 - vybírají věty (je potřeba > 5,000 vět/jazyk)
 - kontrolují věty
 - nahrávají věty (za všechny jazyky je nyní **nahráno 18,000 hodin**, za češtinu **67 hodin**)
 - kontrolují nahrávky vět
- <https://commonvoice.mozilla.org/cs>,
<https://www.mozilla.cz/zpravicky/tag/voice/>

Mozilla Common Voice

- voice.mozilla.org
- iniciativa Mozilly pro vytvoření **komunitního ASR**
- můžete sami přispět k dostupnosti **rozpoznávání řeči** pro váš jazyk
- uživatelé:
 - vybírají věty (je potřeba > 5,000 vět/jazyk)
 - kontrolují věty
 - nahrávají věty (za všechny jazyky je nyní **nahráno 18,000 hodin**, za češtinu **67 hodin**)
 - kontrolují nahrávky vět
- <https://commonvoice.mozilla.org/cs>,
<https://www.mozilla.cz/zpravicky/tag/voice/>

Mozilla Common Voice

- voice.mozilla.org
- iniciativa Mozilly pro vytvoření **komunitního ASR**
- můžete sami přispět k dostupnosti **rozpoznávání řeči** pro váš jazyk
- uživatelé:
 - vybírají věty (je potřeba > 5,000 vět/jazyk)
 - kontrolují věty
 - nahrávají věty (za všechny jazyky je nyní **nahráno 18,000 hodin**, za češtinu **67 hodin**)
 - kontrolují nahrávky vět
- <https://commonvoice.mozilla.org/cs>,
<https://www.mozilla.cz/zpravicky/tag/voice/>

Související technologie

● Dialogové systémy

- Počítačové systémy komunikující s uživatelem pomocí přirozeného jazyka
- Využívají ASR a TTS jako své komponenty

● Rozpoznávání mluvího

- identifikace mluvího – určení, který z registrovaných mluvích pronesl danou větu
- verifikace mluvího – akceptování nebo odmítnutí identity mluvího

● Identifikace mluveného jazyka

- fonémicko-fonetický rozpoznávač pro každý rozpoznávaný jazyk – sledují se fonémy specifické pro každý jazyk
- daná promluva je zpracována všemi rozpoznávači a jako jazyk dané promluvy je zvolen jazyk, jehož rozpoznávač dosáhl nejvyššího skóre

Související technologie

- Dialogové systémy
 - Počítačové systémy komunikující s uživatelem pomocí přirozeného jazyka
 - Využívají ASR a TTS jako své komponenty
- Rozpoznávání mluvího
 - identifikace mluvího – určení, který z registrovaných mluvích pronesl danou větu
 - verifikace mluvího – akceptování nebo odmítnutí identity mluvího
- Identifikace mluveného jazyka
 - fonémicko-fonetický rozpoznávač pro každý rozpoznávaný jazyk – sledují se fonémy specifické pro každý jazyk
 - daná promluva je zpracována všemi rozpoznávači a jako jazyk dané promluvy je zvolen jazyk, jehož rozpoznávač dosáhl nejvyššího skóre

Související technologie

- Dialogové systémy
 - Počítačové systémy komunikující s uživatelem pomocí přirozeného jazyka
 - Využívají ASR a TTS jako své komponenty
- Rozpoznávání mluvího
 - identifikace mluvího – určení, který z registrovaných mluvích pronesl danou větu
 - verifikace mluvího – akceptování nebo odmítnutí identity mluvího
- Identifikace mluveného jazyka
 - fonémicko-fonetický rozpoznávač pro každý rozpoznávaný jazyk – sledují se fonémy specifické pro každý jazyk
 - daná promluva je zpracována všemi rozpoznávači a jako jazyk dané promluvy je zvolen jazyk, jehož rozpoznávač dosáhl nejvyššího skóre

Související technologie

- Dialogové systémy
 - Počítačové systémy komunikující s uživatelem pomocí přirozeného jazyka
 - Využívají ASR a TTS jako své komponenty
- Rozpoznávání mluvčího
 - identifikace mluvčího – určení, který z registrovaných mluvčích pronesl danou větu
 - verifikace mluvčího – akceptování nebo odmítnutí identity mluvčího
- Identifikace mluveného jazyka
 - fonémicko-fonetický rozpoznávač pro každý rozpoznávaný jazyk – sledují se fonémy specifické pro každý jazyk
 - daná promluva je zpracována všemi rozpoznávači a jako jazyk dané promluvy je zvolen jazyk, jehož rozpoznávač dosáhl nejvyššího skóre

Související technologie

- Dialogové systémy
 - Počítačové systémy komunikující s uživatelem pomocí přirozeného jazyka
 - Využívají ASR a TTS jako své komponenty
- Rozpoznávání mluvčího
 - identifikace mluvčího – určení, který z registrovaných mluvčích pronesl danou větu
 - verifikace mluvčího – akceptování nebo odmítnutí identity mluvčího
- Identifikace mluveného jazyka
 - fonémicko-fonetický rozpoznávač pro každý rozpoznávaný jazyk – sledují se fonémy specifické pro každý jazyk
 - daná promluva je zpracována všemi rozpoznávači a jako jazyk dané promluvy je zvolen jazyk, jehož rozpoznávač dosáhl nejvyššího skóre

Související technologie

- Dialogové systémy
 - Počítačové systémy komunikující s uživatelem pomocí přirozeného jazyka
 - Využívají ASR a TTS jako své komponenty
- Rozpoznávání mluvčího
 - identifikace mluvčího – určení, který z registrovaných mluvčích pronesl danou větu
 - verifikace mluvčího – akceptování nebo odmítnutí identity mluvčího
- Identifikace mluveného jazyka
 - fonémicko-fonetický rozpoznávač pro každý rozpoznávaný jazyk – sledují se fonémy specifické pro každý jazyk
 - daná promluva je zpracována všemi rozpoznávači a jako jazyk dané promluvy je zvolen jazyk, jehož rozpoznávač dosáhl nejvyššího skóre

Související technologie

- Dialogové systémy
 - Počítačové systémy komunikující s uživatelem pomocí přirozeného jazyka
 - Využívají ASR a TTS jako své komponenty
- Rozpoznávání mluvího
 - identifikace mluvího – určení, který z registrovaných mluvích pronesl danou větu
 - verifikace mluvího – akceptování nebo odmítnutí identity mluvího
- Identifikace mluveného jazyka
 - fonémicko-fonetický rozpoznávač pro každý rozpoznávaný jazyk – sledují se fonémy specifické pro každý jazyk
 - daná promluva je zpracována všemi rozpoznávači a jako jazyk dané promluvy je zvolen jazyk, jehož rozpoznávač dosáhl nejvyššího skóre

Související technologie

- Dialogové systémy
 - Počítačové systémy komunikující s uživatelem pomocí přirozeného jazyka
 - Využívají ASR a TTS jako své komponenty
- Rozpoznávání mluvího
 - identifikace mluvího – určení, který z registrovaných mluvích pronesl danou větu
 - verifikace mluvího – akceptování nebo odmítnutí identity mluvího
- Identifikace mluveného jazyka
 - fonémicko-fonetický rozpoznávač pro každý rozpoznávaný jazyk – sledují se fonémy specifické pro každý jazyk
 - daná promluva je zpracována všemi rozpoznávači a jako jazyk dané promluvy je zvolen jazyk, jehož rozpoznávač dosáhl nejvyššího skóre

Související technologie

- Dialogové systémy
 - Počítačové systémy komunikující s uživatelem pomocí přirozeného jazyka
 - Využívají ASR a TTS jako své komponenty
- Rozpoznávání mluvího
 - identifikace mluvího – určení, který z registrovaných mluvích pronesl danou větu
 - verifikace mluvího – akceptování nebo odmítnutí identity mluvího
- Identifikace mluveného jazyka
 - fonémicko-fonetický rozpoznávač pro každý rozpoznávaný jazyk – sledují se fonémy specifické pro každý jazyk
 - daná promluva je zpracována všemi rozpoznávači a jako jazyk dané promluvy je zvolen jazyk, jehož rozpoznávač dosáhl nejvyššího skóre