

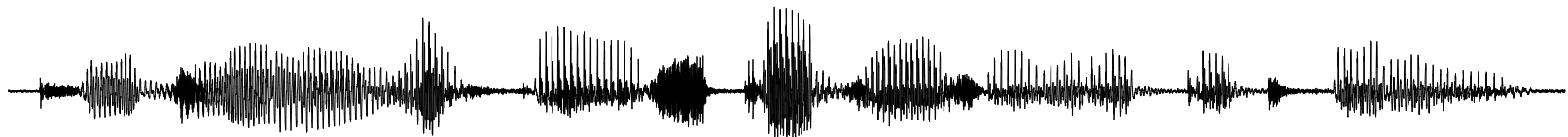
Grafeem-foneemomzetting voor spraaksynthese

Two steps

- P G & E will file schedules on April 20.
- **TEXT ANALYSIS:** Text into intermediate representation:

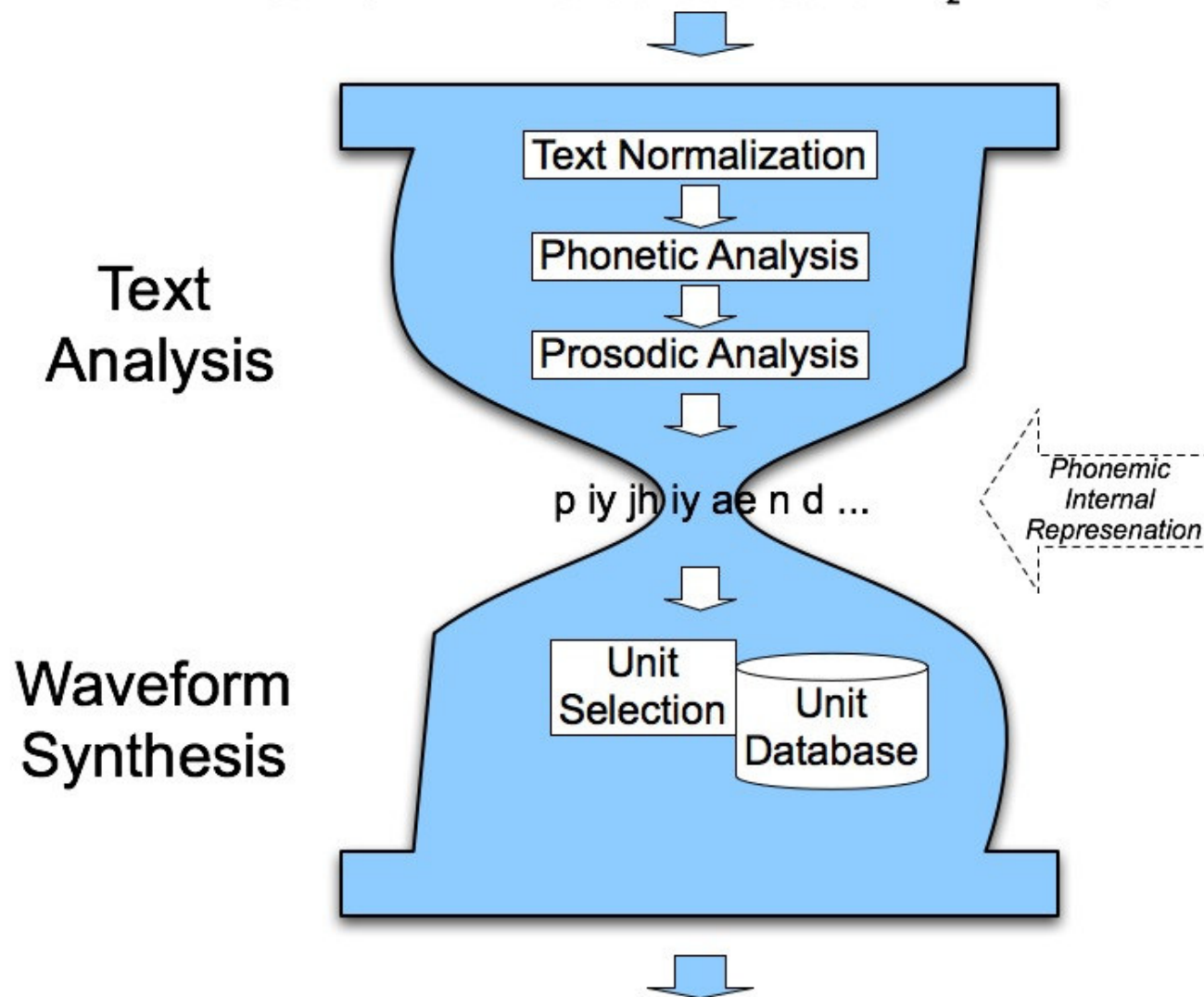
| | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
|---|----|-----|----|------|------|---|----|-------|----|------|---|----|---|---|---|----|----|----|---|---|----|---|----|---|---|----|---|---|---|----|---|---|----|----|----|
| P | G | AND | * | WILL | FILE | * | ON | APRIL | * | L-L% | | | | | | | | | | | | | | | | | | | | | | | | | |
| p | iy | jh | iy | ae | n | d | iy | w | ih | l | f | ay | l | s | k | eh | jh | ax | l | z | aa | n | ey | p | r | ih | l | t | w | eh | n | t | iy | ax | th |

- **WAVEFORM SYNTHESIS:** From the intermediate representation into waveform



The Hourglass

PG&E will file schedules on April 20.



fasen in tekst-naar-spraak

- **Tekstbewerking**
 - niet-letters, afkortingen, niet-Nederlandse woorden
- **Morfologische decompositie**
 - samenstellingen, syllabegrenzen [morfologie]
- **Grafeem-foneemomzetting**
 - lettersymbolen -> klanksymbolen [fonologie]
- **Melodie en ritme (prosodie)**
 - woordklemtoon [fonologie] en zinsaccenten [parsing, semantiek]
- **Synthese**
 - akoestische realisatie [fonetiek]

gebruikelijk ontwikkeltraject

- Dag: eerste 50%
- Week: volgende 25 %
- Maand: volgende 10 %
- Jaar: volgende 10 %
- ?: laatste 5%

tekstbewerking: cijfers

- **getallen,**

5 -> "vijf", 21 -> "een en *twintig*" (twenty-one)

- **telefoonnummers**

020-6680470 -> "nul *twintig zes zes tachtig vier zeventig*"

- **geldbedragen**

euro 32.54 -> "twee en dertig *euro vier en vijftig*"

- **huisnummers**

13-hs -> "dertien huis"

- **Romeinse cijfers**

MDCCLXIV -> "zeventien honderd vier en zestig"

tekstbewerking: afkortingen

- als losse letters

KLM -> "K L M"

- als woord

VARA -> "vara"

- als afgekort woord

tel. -> "telefoon"

maar let op de juiste interpretatie van de punt
(telefoon vs tel – *de laatste tel.*)

tekstbewerking: leestekens

- punten, komma's, puntkomma, etc -> prosodie

foneem-grafeemomzetting

symbolen voor spraakklanken

- fonetisch alfabet

IPA International Phonetic Alphabet

- lastige tekenset in pre-Unicode tijdperk

œ ø ʌ ʒ tʃ ʒ ɸ ...

- computer fonetisch alfabet

verschillende voorbeelden, bv SAMPA

(speech assessment methods phonetic alphabet)

spreek uit en luister naar je klanken!

3. Letter-to-Sound: Getting from words to phones

- Two methods:
 - Dictionary-based
 - Rule-based (Letter-to-sound=LTS)
- Early systems, all LTS
- MITalk was radical in having huge 10K word dictionary
- Now systems use a combination

Pronunciation Dictionaries: CMU

- CMU dictionary: 127K words
 - <http://www.speech.cs.cmu.edu/cgi-bin/cmudict>
- Some problems:
 - Has errors
 - Only American pronunciations
 - No syllable boundaries
 - Doesn't tell us which pronunciation to use for which homophones
 - (no POS tags)
 - Doesn't distinguish case
 - The word US has 2 pronunciations
 - [AH1 S] and [Y UW1 EH1 S]

Pronunciation Dictionaries: UNISYN

- UNISYN dictionary: 110K words (Fitt 2002)
 - <http://www.cstr.ed.ac.uk/projects/unisyn/>
 - going: { g * ou }.> i ng >
 - antecedents: { * a n . t ^ i . s ~ ii . d n ! t } > s >
 - dictionary: { d * i k . sh @ . n ~ e . r ii }
- Benefits:
 - Has syllabification, stress, some morphological boundaries
 - Pronunciations can be read off in
 - General American
 - RP British
 - Australia
 - Etc
 - (Other dictionaries like CELEX not used because too small, British-only)

Dictionaries aren't sufficient

- Unknown words (= OOV = “out of vocabulary”)
 - Increase with the (sqrt of) number of words in unseen text
 - Black et al (1998) OALD on 1st section of Penn Treebank:
 - Out of 39923 word tokens,
 - 1775 tokens were OOV: 4.6% (943 unique types):

| names | unknown | Typos/other |
|-------|---------|-------------|
| 1360 | 351 | 64 |
| 76.6% | 19.8% | 3.6% |

- So commercial systems have 4-part system:
 - Big **dictionary**
 - **Names** handled by special routines
 - **Acronyms** handled by special routines (previous lecture)
 - Machine learned **g2p** algorithm for other unknown words

Names

- Big problem area is names
- Names are common
 - 20% of tokens in typical newswire text will be names
 - 1987 Donnelly list (72 million households) contains about 1.5 million names
 - Personal names: McArthur, D'Angelo, Jiminez, Rajan, Raghavan, Sondhi, Xu, Hsu, Zhang, Chang, Nguyen
 - Company/Brand names: Infinit, Kmart, Cytoc, Medamicus, Inforte, Aeon, Idexx Labs, Bebe

Names

- **Methods:**
 - Can do morphology (Walters -> Walter, Lucasville)
 - Can write stress-shifting rules (Jordan -> Jordanian)
 - Rhyme analogy: Plotsky by analogy with Trostsky (replace tr with pl)
 - Liberman and Church: for 250K most common names, got 212K (85%) from these modified-dictionary methods, used LTS for rest.
 - Can do automatic country detection (from letter trigrams) and then do country-specific rules
 - Can train **g2p** system specifically on names
 - Or specifically on types of names (brand names, Russian names, etc)

Letter-to-Sound Rules

- Earliest algorithms: handwritten Chomsky+Halle-style rules:
 - $c \rightarrow [k] / \text{---} \{a,o\}V$; context dependent
 - $c \rightarrow [s]$; context independent
- Festival version of such LTS rules:
 - (LEFTCONTEXT [ITEMS] RIGHTCONTEXT = NEWITEMS)
- Example:
 - (# [c h] C = k)
 - (# [c h] = ch)
- # denotes beginning of word
- C means all consonants
- Rules apply in order
 - “christmas” pronounced with [k]
 - But word with ch followed by non-consonant pronounced [ch]
 - E.g., “choice”

symbool omzetting

- 1 naar 1
p -> p , a -> A
- 1 naar 2
o -> O~
- 2 naar 1
(*dus onderzoek altijd eerst context!*)
ch -> x , sj -> S , ee -> e , ie -> i , ng -> N
- 2 naar 2
ij -> Ei

contekst afhankelijkheid

- positie in woord

omringende letters

ban -> bAn bang -> bAN

final devoicing

dak -> dAk bad -> bAt

woordbegin <-> woordeinde

lat -> lAt tal -> tAL

wit -> wlt ruw -> ryW

assimilatie

zakdoek -> zAgduk en geen zAkduk

Voorbeelden

C*V*C* woorden zonder diacrieten

- positie van letter in woord
- positie van letter in cluster
- onset-nucleus-coda identificatie helpt
 - Klinkerclusters
 - l-eeu-w, g-eëe-rd
 - Specifieke medeklinkerclusters
 - Fonotactische restricties ->
ook voor syllabegrenzen, maar vgl wegrennen

Contextgevoeligheid (C) *in* $C^*V^*C^*$

- Ongevoelig *in monosyllabe*
 - f, k, m, p, r, t, v, x, z
- Volgende context
 - b (#), c (varia), d (#), g (**goa**), j (#), l (#),
n (**ng**, #), q (**qu**), s (**sj**), w (#)
- Vorige context
 - h (**ch**), j (**ij**), g (**ng**) **vet als segment nemen**

Contextgevoeligheid (V) in $C^*V^*C^*$

Veel op te lossen door segmentdefinitie

- aai, eeu, ooi, ieu
 - mooi-er , ui-en
- aa, ee, oo, uu, ie, oe, eu [1 symbool]
ai, oi, au, ou, ij, ei, ui [2 symbolen]
ua, io [grens $uV^* iV^*$]
- korte klinkers a, i, o, u, y
afhandelen na lange klinkers en tweeklanken
 - kan, act = uitzondering
- e
 - lastig: e, E, @, l

Regelvolgorde

mooier

- eerst langste cluster: ooi

m-ooi-e-r -> moj@r (mO:j@r)

- en niet

m-oo-ie-r -> moir

of

m-o-o-i-e-r -> mOOIEr

| | plofklanken | wrijfklanken | nasalen |
|---------------------------------|--|---|---------------------|
| labiaal (lippen) | p (pad) b (bad) | f (fiets) v (vat) | m (mat) |
| alveolair (achter de tanden) | t (tak) d (dak) | s (sap) z (zat) | n (nat) |
| palataal (verhemelte) | tj (potje) dj (djintan) | S (sjaal) Z (plantage) | nj (anjer) |
| velair (achter) | k (kat) g (zakdoek, goal) | x (lach) G (gat) | N (bang) |
| glottaal (stem) | _ (stilte) | h (huis) | |

Liquida en halfvokalen

| | woordbegin | woordeinde |
|-------------|---------------------------|---|
| Liquida | l (lang, alle) r (rug) | L (al) r (haar) R (haar) |
| Halfvokalen | w (wit) j (jan) | W (sneeuw) na / i, e, y / J (aai) |

Korte klinkers

| | voor | | achter | |
|--------|---------|---------|----------|---------|
| | [-rond] | [+rond] | [-rond] | [+rond] |
| hoog | I (bid) | Y (put) | | |
| midden | E (bed) | | @ (rede) | O (bos) |
| laag | | | A (bak) | |

Lange klinkers

| | voor | | achter | |
|--------|----------|----------|----------|----------|
| | [-rond] | [+rond] | [-rond] | [+rond] |
| hoog | i (bied) | y (buut) | | u (boek) |
| midden | e (beet) | 2 (beuk) | | o (boot) |
| laag | | | a (baat) | |

Tweeklanken

| | voor | | achter | |
|---------------|----------------|----------------|----------------|----------------|
| | [-rond] | [+rond] | [-rond] | [+rond] |
| midden | Ei (bijt) | 9y (buit) | | Au (bout) |

Gekleurde klinkers voor r

| | voor | | achter | |
|---------------|------------------|------------------|----------------|------------------|
| | [-rond] | [+rond] | [-rond] | [+rond] |
| midden | I: (beer) | Y: (deur) | | O: (door) |

Franse klanken

| | voor | | achter | |
|--|-------------------|----------------|-----------------|-----------------|
| | [-rond] | [+rond] | [-rond] | [+rond] |
| | E~ (mannequin) | Y~ (parfum) | A~ (chanson) | O~ (chanson) |

[Oe] uit freule