



Multimodal speech synthesis/ (N)GSLT - Speech Technology course
Björn Granström
CTT, KTH

School for Computer Science and Communication




KTH - Kungliga tekniska högskolan

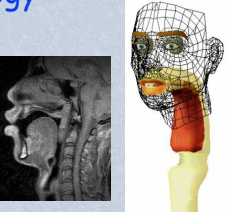
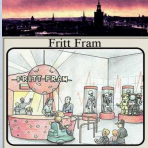
Department of Speech, Music and Hearing

Multimodal speech synthesis - (NGSLT 2006 [2])

CTT - Centre for speech technology

Research areas

- Speech production
- Speech perception
- Communication aids
- Multimodal speech synthesis
- Speech understanding
- Speaker characteristics
- Interactive dialog systems
- Affective computing
- Language learning

Multimodal speech synthesis - (NGSLT 2006 [3])

CTT - modes of operation

- long-term research projects,
- participation by CTT personnel in international projects,
- research exchange with leading foreign groups,
- graduate-level research training,
- dissemination of competence through the exchange of researchers within Sweden and collaboration with other Swedish research groups
- special information dissemination efforts

Multimodal speech synthesis - (NGSLT 2006 [4])

CTT partners phase 4 - 2004-2006

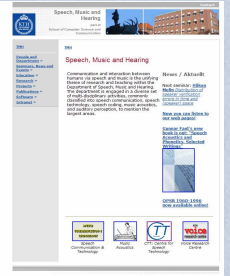


Babel-Infovox
 English Town
 GN Resound
 Hjälpmedelsinstitutet
 HoneySoft
 IcePeak
 Dolphin Audio Publishing
 LingTek
 Phoneticom
 Polycom Technologies
 SaabTech
 SpeechCraft
 STTS
 Sveriges Radio
 Sveriges Television
 TeliaSonera
 TPB - Talboks- och punktskriftsbiblioteket
 Vattenfall
 VoiceProvider

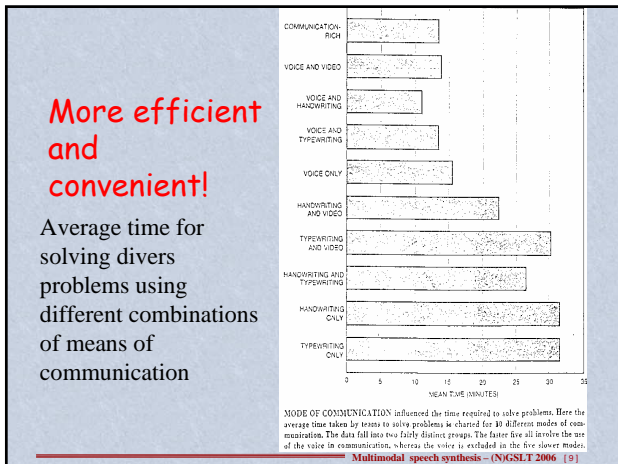
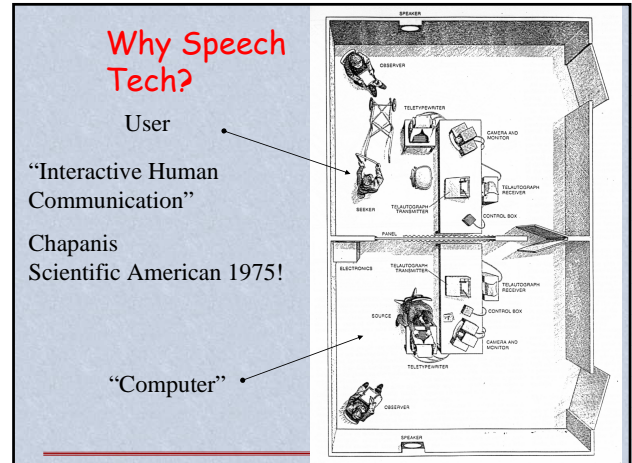
Multimodal speech synthesis - (NGSLT 2006 [5])

The work presented in this lecture is the result of many researchers' efforts at the Department of Speech, Music and Hearing

Reports and further information can be found on our home page www.speech.kth.se



Multimodal speech synthesis - (NGSLT 2006 [6])



Speech is fast and verbose

	COMMUNICATION-RICH	VOICE	HANDWRITING	TYPEWRITING	
				EXPERIENCED TYPISTS	INEXPERIENCED TYPISTS
SOLUTION TIME IN MINUTES	29	33	53.3	66.2	69
NUMBER OF MESSAGES	230.4	163.8	15.9	27.2	31.5
NUMBER OF SENTENCES	372.6	275.9	24.9	45.8	44.1
TOTAL NUMBER OF WORDS	1,583.8	1,374.8	224.8	322.9	287.4
TOTAL NUMBER OF DIFFERENT WORDS	397.5	305.9	118.5	150.5	133.4
TYPE-TOKEN RATIO	.3	.3	.6	.5	.6
NUMBER OF WORDS PER MINUTE	190.3	171.2	17.3	18.1	10.2

EXPERIMENTAL RESULTS are enumerated for the solution of problems by various modes of communication. "Type-token ratio" is ratio of different words to total words. Problem solving by voice takes the least time but is wordier than the other modes are.

Multimodal speech synthesis - (NGSLT 2006 [10])

Speech technology is making money!

Classic example : AT&T and Lucent Technologies VRCP, "Voice Recognition Call Processing" service

- Selection of payment method
- Vocabulary only five words : collect, calling card, third number, person, operator
- More than 5 000 000 calls per day
- Earnings already (1999?) more than AT&T/Bell labs' total investment in speech research

Multimodal speech synthesis - (NGSLT 2006 [11])



Speech synthesis developments at KTH

Multimodal speech synthesis - (NGSLT 2006 [13])

The KTH speech group - Early days

Gunnar Fant and OVE I 1953

Multimodal speech synthesis - (NGSLT 2006 [14])

OVE I, in the WaveSurfer tool

- Interface is based around *WaveSurfer*, a general purpose tool for speech and audio viewing, editing and labelling
- TTS and Talking Head functionality is added as plug-ins
- WaveSurfer (presently without TTS&TH) works on all common platforms and is freely available as open source
- Modules from Waves available – formants and F0 in present release – thanks to Microsoft and AT&T

<http://www.speech.kth.se/wavesurfer>

Multimodal speech synthesis - (NGSLT 2006 [15])

Ove II, 1958

1961

1962

Multimodal speech synthesis - (NGSLT 2006 [16])

KTH/TTS history

- 1967, Digitally controlled OVE III
- 1974, Rule-based system RULSYS
 - transformation rules
- 1979, Mobile text-to-speech system
 - used by a non-vocal child
- 1982, Portable TTS (ICASSP, Paris)
 - Multilingual
 - MC 68000, NEC 7720
- 1983, Founding of Infovox Inc.

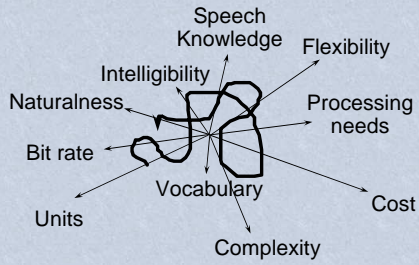
Multimodal speech synthesis - (NGSLT 2006 [17])

What do we mean by Speech synthesis?

- **Recorded speech**
 - Words or phrases (*telephone banking*)
 - Fixed vocabulary – maintenance problems...
- **Concatenative speech synthesis**
 - Diphones or larger units (unit selection)
 - LPC: source filter model (too simple?)
 - PSOLA/MBROLA/HNM – and rules for prosody
 - One speaker
- **Parametric synthesis**
 - Formant synthesis
 - Articulatory synthesis
 - flexible
 - But lower quality - today
- **Multimodal synthesis**

Multimodal speech synthesis - (NGSLT 2006 [18])

The synthesis space



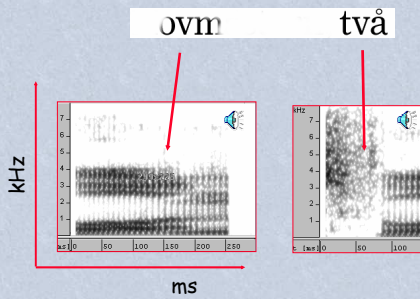
Multimodal speech synthesis - (NGSLT 2006 [19])

TEXT vs SPEECH

- Parallel vs. sequential
- permanent vs disappearing
- Text as transcription of speech?
- Example - WaveSurfer demo :
 - Palindrome
 - "Var sak har två sidor"

Multimodal speech synthesis - (NGSLT 2006 [20])

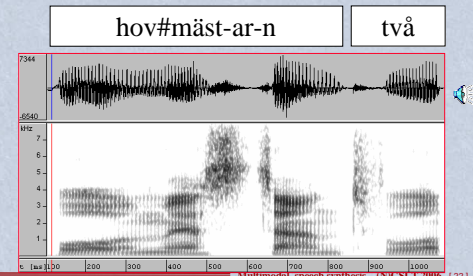
Phonetics



Multimodal speech synthesis - (NGSLT 2006 [21])

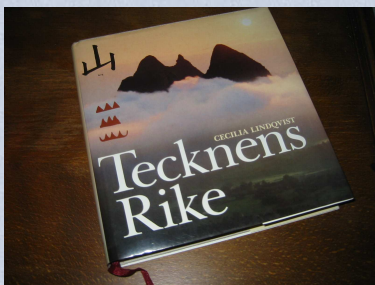
Phonology, Morphology

Hovmästarn, två



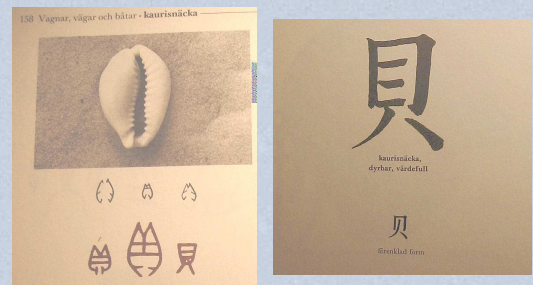
Multimodal speech synthesis - (NGSLT 2006 [22])

Non-phonetic writing



Multimodal speech synthesis - (NGSLT 2006 [23])

Signs vs. technology



Multimodal speech synthesis - (NGSLT 2006 [24])

Is the sign what you think?

Människan • man 25

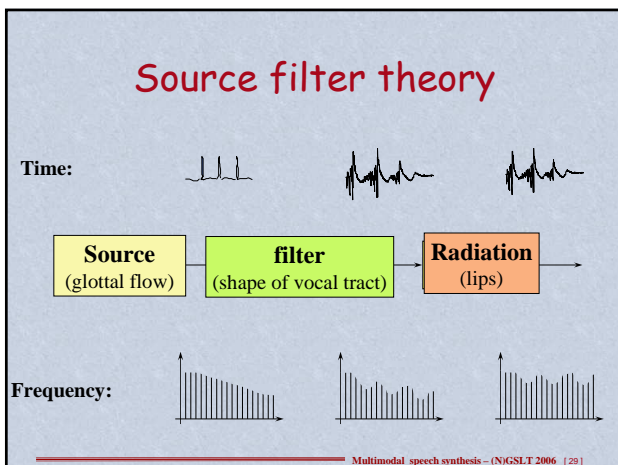
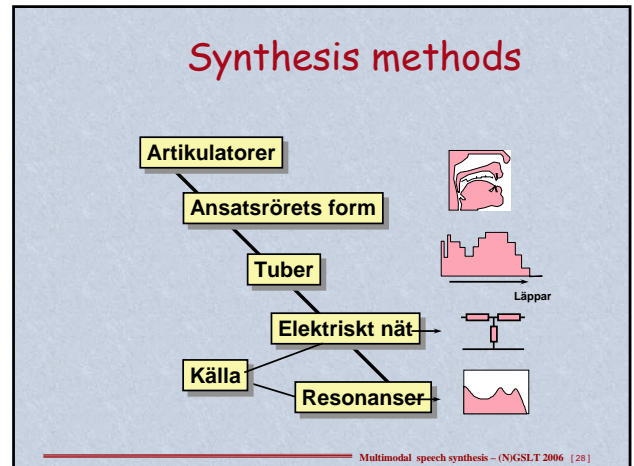
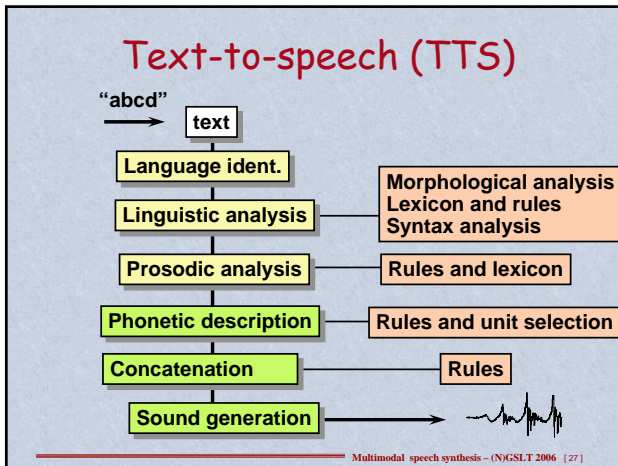
man, människa

Multimodal speech synthesis - (NGSLT 2006 [25])

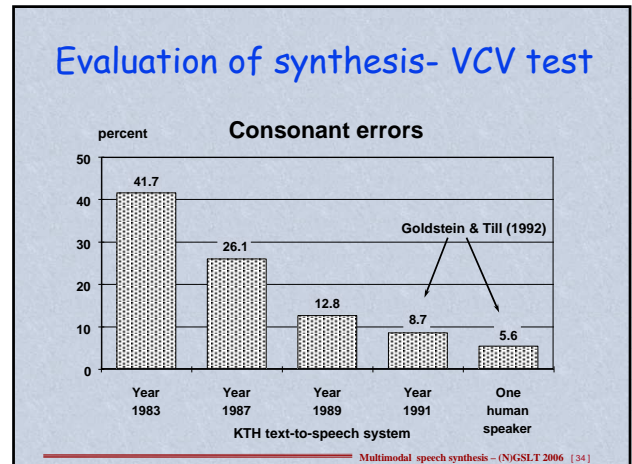
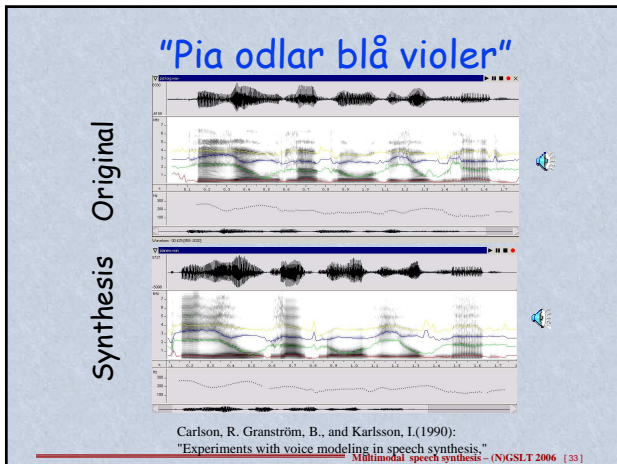
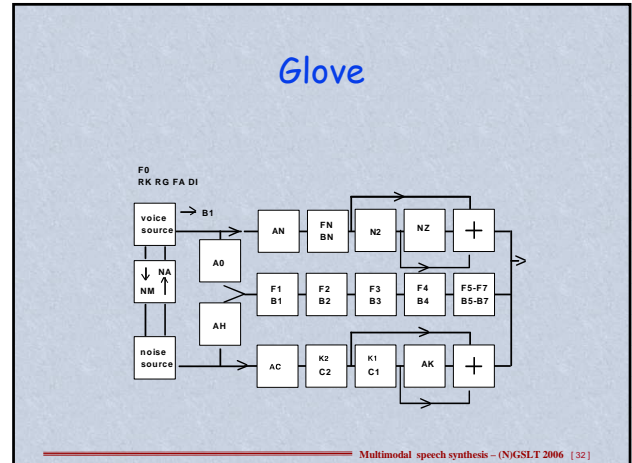
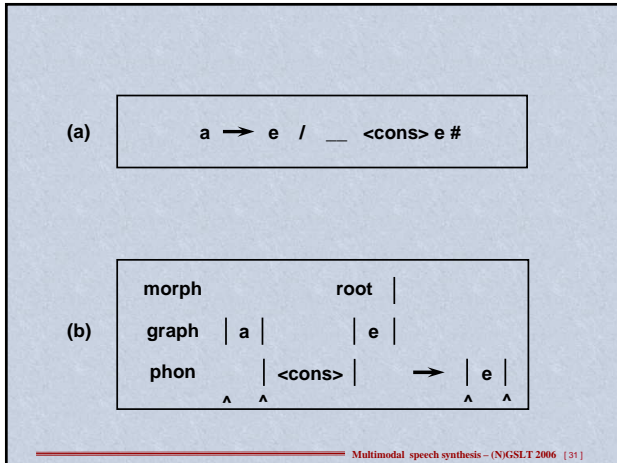
How is a Chinese lexicon organized?

man, människa	人 / 人
stor	大 / 大
öga	目 / 目
ansikte, yta	面 / 面
öra	耳 / 耳
näsa, själv	自 / 自
mun	口 / 口
tand, tänder	齒 / 齒

Multimodal speech synthesis - (NGSLT 2006 [26])



- ### RULSYS Rules - Features
- Is IPA synthesiser possible?
 - Based on generative phonology
 - Language specific definitions
 - Rules for contextual modifications
 - Examples
 - Interactive rule manipulations
- Multimodal speech synthesis - (NGSLT 2006 [30])



Speaker characteristics

Speaker
dialect, sex, social, education, age

Situation
formality, style, interspeaker relation

Complex description
many dependent variables

Multimodal speech synthesis - (NGSLT 2006 [35])

- ## Speaker characteristics
- TIDE/Vaess project (Voices, attitudes and emotions in synthetic speech)
 - "Voice fitting"
 - Software, Glove synthesis
 - 10 user controlled parameters
 - Rule specified connection to synthesis parameters
- Multimodal speech synthesis - (NGSLT 2006 [36])





Emotional/ubiquitous computing - do we want it?



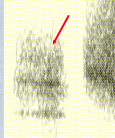


Early BBC vision - the conversational toaster
Thanks to Mark Huckvale, UCL, for the video clip

Multimodal speech synthesis - (NGSLT 2006 [37])

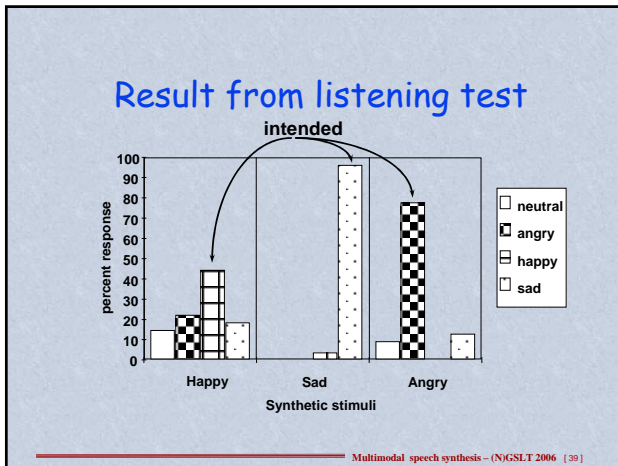
Emotions

natural synthesis   Neutral  natural synthesis 


Happy

Angry   natural synthesis 

Multimodal speech synthesis - (NGSLT 2006 [38])



Emotions

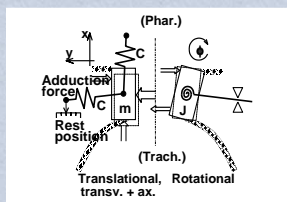


Happiness Anger Fear

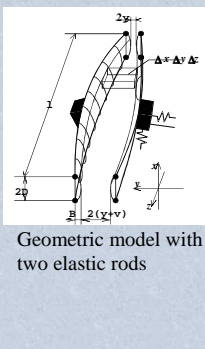
Surprise Disgust Sadness

Multimodal speech synthesis - (NGSLT 2006 [40])

Voice source simulation



Mechanical model



Geometric model with two elastic rods

Air flow through glottis

Multimodal speech synthesis - (NGSLT 2006 [41])

Voice source simulation examples

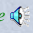
Samplingsintervall: 1 parameter


Stämbandsanatomi: 18 parametar, t ex längd, vikt, dämpningsfaktor


Stämbandsartikulation: 7 parametar, t ex restglapp, lungtryck


Högtrycksartikulation: 4 parametar, t ex tonhöjd, ljudstyrka

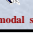
Talrörsartikulation: 8 parametar, t ex area, dämpningsfaktor


Exempel *originallyttrande* 

fördubblad stämbandsmassa 

förlängning av stämband 13->17 mm 

assymetriförändring 1.0->1.8 

assymetriförändring 1.0->2.0 

restgap 0.1->1.0 mm 

Multimodal speech synthesis - (NGSLT 2006 [42])

Predictable word accents

- Word accent and stress
 - 'anden (bird)
 - 'änden (spirit)
 - 'änk-pasta
 - 'änk-pastej
 - 'änk-pastejen
 - 'änk-pastejs-mål-tid
- Name pronunciation
 - karl-'èrik (karl-'erik) + 'anne-marie
 - Onomastica/EU
- Reduction
 - be'töng-väg(g)s-konstruktion (C-C Elert)

Multimodal speech synthesis - (NGSLT 2006 [43])

Letter-to-sound vs. lexicon

Size vs. Precision

Maintenance - “half of the words are unique”

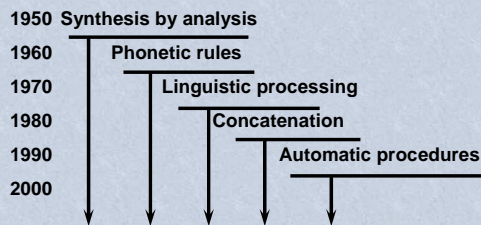
Name pronunciation (Onomastica Project)

Misspellings

Use rules, morphology, analogy...as fall back

Multimodal speech synthesis - (NGSLT 2006 [44])

Research trends in speech synthesis



Multimodal speech synthesis - (NGSLT 2006 [45])

Concatenative synthesis

- Already Peterson et al. (1958)
- Dixon and Maxey (1968)
- “Diacic Units”, (Olive, 1977)
- “PSOLA”, Charpentier and Stella (1986)
- Review, Möbius (2000)
- Unit selection
 - Concatenation cost
 - Quality/size of database
 - Prosody/speaking styles

Multimodal speech synthesis - (NGSLT 2006 [46])

Concatenative synthesis Signal manipulations

- Prosodic modifications
 - Possibility to modify F0
 - Possibility to lengthen or shorten segments
- Spectral modifications
 - Interpolation of spectrum at joints
- Early technique - LPC

Multimodal speech synthesis - (NGSLT 2006 [47])

Speak&Spell - TexasInstrument Christmas 1978



Multimodal speech synthesis - (NGSLT 2006 [48])

PSOLA



- Pitch pulses moved in time to fit F0 contour
- Conceptually simple and computationally efficient
 - Need for precise pitch pulse marking
 - Could not handle spectral interpolation

Multimodal speech synthesis - (NGSLT 2006 [49])

Unit selection

- Large databases of recorded natural speech
- Minimal processing
- Annotation of database – what information is needed?
- Synthesis defaults to transcription and search problem
- Few cuts > maximally long units selected (but context and prosody must fit well)
- Target and concatenation costs

Multimodal speech synthesis - (NGSLT 2006 [50])

Synthesis methods

- **Unit selection – minimal processing**
 - chatr(ATR) weather (CSTR) good, less good
- **Diphone synthesis**
 - Svensk(Mbrola), Fransk (ICP) (CNET)
- **Formantsynthesis**
 - Svensk (KTH), Fransk (KTH)

Multimodal speech synthesis - (NGSLT 2006 [51])

Unit selection - BrightSpeech

- Swedish
- Norwegian

Multimodal speech synthesis - (NGSLT 2006 [52])

Examples of Synthesized Speech

Universität Stuttgart, Institut für Maschinelle Sprachverarbeitung

[German] [English] [French] [Dutch] [Spanish] [Italian] [Portuguese]
[Swedish] [Norwegian] [Finnish] [Estonian] [Icelandic] [Czech]
[Russian] [Greek] [Croatian] [Romanian] [Japanese] [Chinese]
[Korean] [Hebrew] [Arabic]

<http://www.ims.unistuttgart.de/~moehler/synthespeech/examples.html>

also e.g. *Synthesis examples:*

<http://www.naturalvoices.att.com/>

<http://www.acapela-group.com/demos/demos.asp>

<http://www.naturalvoices.att.com/>

<http://www.nextup.com/>

Multimodal speech synthesis - (NGSLT 2006 [53])

To make it work today: Hybrid systems
e.g. Who has number nn....?
Key input - Synthesis mixed with
recorded speech (Call +46 118 999)



Multimodal speech synthesis - (NGSLT 2006 [54])

Hybrid methods, cont.

- Rolf Carlson, Tor Sigvardson, Arvid Sjölander (2002). Data-driven formant synthesis. *Proc of Fonetik 2002, TMH-QPSR*
- David Öhlin and Rolf Carlson Data-driven Formant Synthesis, *Proc of Fonetik 2004*
- + Four MSc theses (Sigvardson, Sjölander, Vinet, Öhlin)
- All available on www.speech.kth.se

Multimodal speech synthesis - (NGSLT 2006 [55])

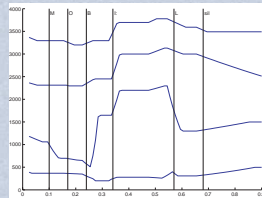
Aim

- Keeps the flexibility of the formant synthesis
- More natural sounding than rule-driven synthesis
- Speaker adaption

Multimodal speech synthesis - (NGSLT 2006 [56])

Rule-driven formant synthesis

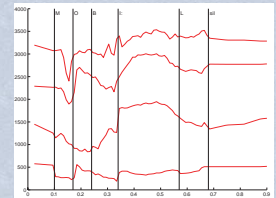
- Parameters are generated by rule (RULSYS, Carlson et al.)
- Formant values are generated by interpolating between target frequencies
- Parameters are fed to a synthesizer (GLOVE, Carlson et al.)



Multimodal speech synthesis - (NGSLT 2006 [57])

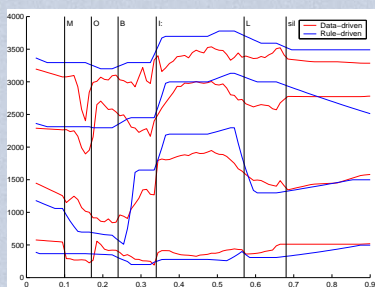
Data-driven formant synthesis

- Some parameters (namely, the first four formants) are replaced by data
- Same synthesizer



Multimodal speech synthesis - (NGSLT 2006 [58])

Synthesis comparison



Multimodal speech synthesis - (NGSLT 2006 [59])

Data-driven formant synthesis

- Formants are replaced through unit selection from a formant diphone library
- Formant trajectories are scaled and interpolated to fit the rule-generated durations

Multimodal speech synthesis - (NGSLT 2006 [60])

Cost function

- Designed to promote probable formant candidates
- Penalizes:
 - Large bandwidths
 - Large frequency deviations (given the current phoneme)
 - Large frequency jumps (promotes smooth trajectories)

Multimodal speech synthesis - (NGSLT 2006 [61])

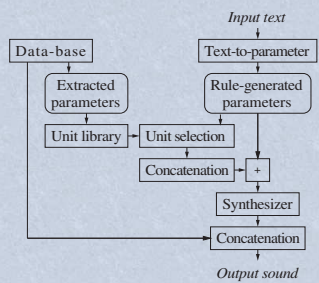
Voiceless consonants

- Replace the voiceless fricatives and plosives with recorded versions
- Voiceless fricatives in Standard Swedish: /f/, /s/, /sj/, /tj/, and /rs/
- Voiceless plosives: /k/, /p/, /t/, /t/

(/h/ is excluded)

Multimodal speech synthesis - (NGSLT 2006 [62])

Text-to-speech synthesis



Multimodal speech synthesis - (NGSLT 2006 [63])

Listening test evaluation 1

- 15 subjects, 20 sentences, continuous scale
- Data-driven synthesis with non-corrected formant data was judged more natural sounding than rule-driven synthesis




Multimodal speech synthesis - (NGSLT 2006 [64])



Listening test evaluation 2

- 12 subjects, 10 sentences, binary scale
- Data-driven synthesis with manually corrected formant data was preferred in 73 % of the cases over rule-driven synthesis

Multimodal speech synthesis - (NGSLT 2006 [65])

Sound samples

- Strindberg, rule-driven: 
- Strindberg, data-driven: 
- Strindberg, MBROLA: 

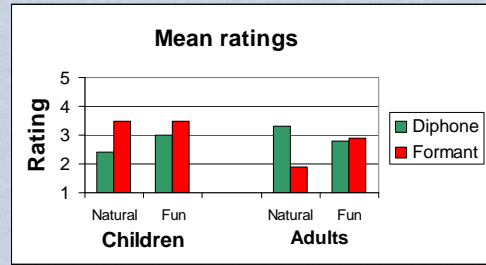
- "Pytteliten", rule-driven: 
- "Pytteliten", data-driven: 

Multimodal speech synthesis - (NGSLT 2006 [66])

Child-directed speech synthesis

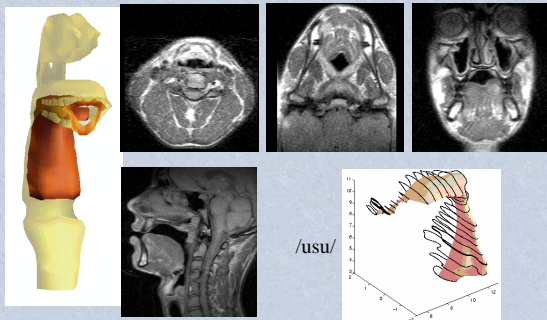
- Increase prosodic variation in synthesis
- How do children between the ages of 9 and 11 react to:
 - default, F0, duration?
 - diphone synthesis vs formant synthesis?

Multimodal speech synthesis - (NGSLT 2006 [67])



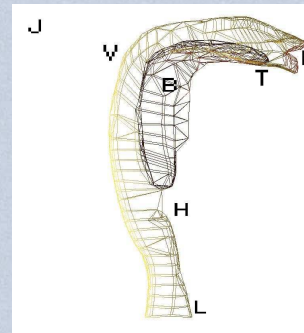
Multimodal speech synthesis - (NGSLT 2006 [68])

Articulatory synthesis



Multimodal speech synthesis - (NGSLT 2006 [69])

Articulatory parameters



- Jaw opening
- Lip rounding
- Lip Protrusion
- Tongue position
- Tongue height
- Tongue tip
- Velum
- Hyoid

Multimodal speech synthesis - (NGSLT 2006 [70])

Articulatory synthesis potential use



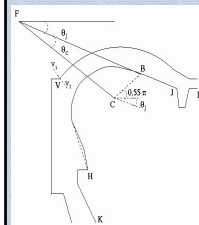
- Articulatory synthesis
 - Calculations directly from cross sectional areas
 - Fluid dynamics calculations



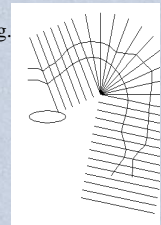
- Visual synthesis
 - Articulation training
- Demonstrations and research

Multimodal speech synthesis - (NGSLT 2006 [71])

Why a 3D model?



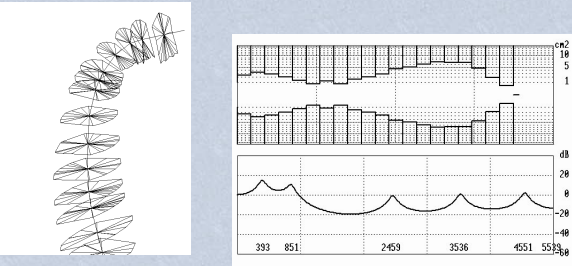
- 2D articulatory models, e. g.
 - Mermelstein
 - Maeda
- Imposing a third dimension:
 - area = $a \cdot (\text{width})^b$
- A 3D model:
 - Direct calculation
 - laterals etc.



Multimodal speech synthesis - (NGSLT 2006 [72])

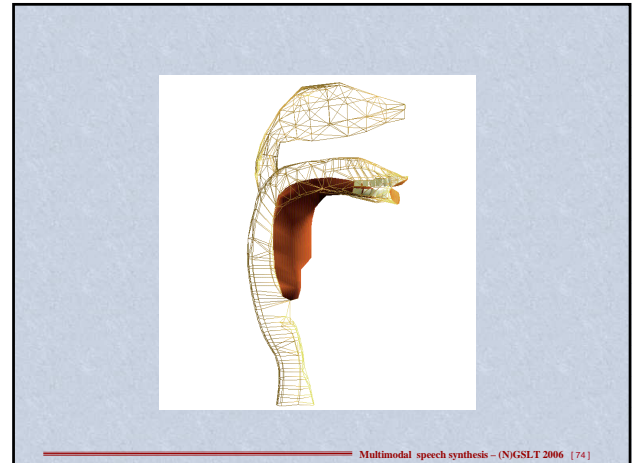
From areas to formants

A transfer function is determined from the cross-sectional areas

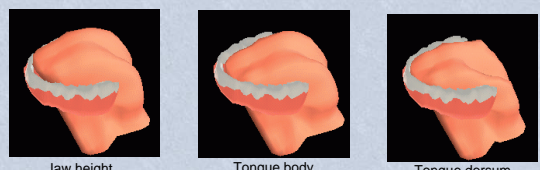


The diagram illustrates the relationship between the physical cross-sectional areas of the vocal tract and the resulting acoustic properties. On the left, a series of wireframe cross-sections shows the shape of the vocal tract. On the right, three plots are shown: the top plot is a bar chart of cross-sectional areas, the middle plot is a bar chart of reflection coefficients, and the bottom plot is a line graph of the transfer function showing formant frequencies. The x-axis for the bottom plot is labeled with values 393, 851, 2459, 3536, 4551, and 5538.

Multimodal speech synthesis - (NGSLT 2006 [73])



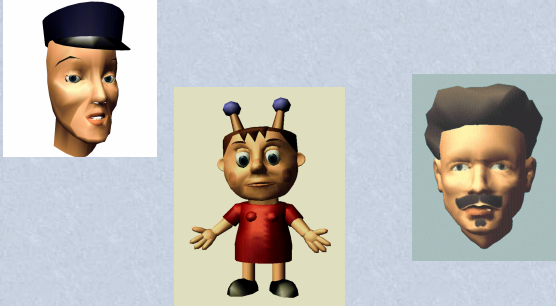
Ranges of parameter activations



Three 3D models of the vocal tract, showing different parameter activations. The models are rendered in a light blue color and are shown in a profile view, facing right. The first model is labeled "Jaw height", the second is labeled "Tongue body", and the third is labeled "Tongue dorsum".

Multimodal speech synthesis - (NGSLT 2006 [75])

Multimodal speech synthesis

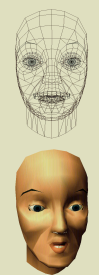


Three 3D models of human faces, showing different styles. The first model is a realistic face with a black cap. The second model is a cartoon character with a red shirt and black pants. The third model is a stylized face with a black cap and a mustache.

Multimodal speech synthesis - (NGSLT 2006 [76])

Talking heads - Applications

- Improved speech synthesis
- Human-Computer Interface in spoken dialogue systems
- Aid for hearing impaired
- Educational software
- Stimuli for perceptual experiments
- Entertainment: games, virtual reality, movies etc.




Two 3D models of human faces. The top model is a wireframe face, and the bottom model is a realistic face.

Multimodal speech synthesis - (NGSLT 2006 [77])

A new paradigm for human-computer interaction

- Shift from desktop-metaphor to person-metaphor
- Spoken dialogue as well as non-verbal communication
- Take advantage of the user's social skills
- Strive for believability, but not necessarily realism



3D model of a human face with the text "Listening & Thinking" below it.

Multimodal speech synthesis - (NGSLT 2006 [78])

Conventions? - use same as for person-to-person communication



SONY SDR



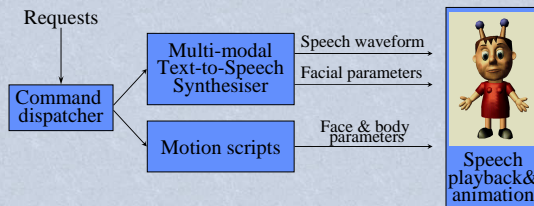
Multimodal

Tasks of an Animated Agent

- Provide intelligible synthetic speech
- Indicate emphasis and focus in utterances
- Support turn-taking
- Give spatial references (gaze, pointing etc)
- Provide non-verbal back-channeling
- Indicate the system's internal state

Multimodal speech synthesis - (NGSLT 2006 [80])

Animated Character - architecture



Multimodal speech synthesis - (NGSLT 2006 [81])

Parameters used for articulatory control of the face.

- Jaw rotation
- Lip rounding
- Lip protrusion
- Mouth width
- Bilabial closure
- Labiodental closure
- Upper lip raise
- Lower lip depression
- Apex
- Tongue length
- + more for prosody, attitude, emotions, turn-taking, back-channeling, pointing

Multimodal speech synthesis - (NGSLT 2006 [82])

The WaveSurfer Tool



- Interface is based around *WaveSurfer*, a general purpose tool for speech and audio viewing, editing and labelling
- TTS and Talking Head functionality is added as plug-ins
- WaveSurfer (presently without TTS&TH) works on all common platforms and is freely available as open source

<http://www.speech.kth.se/wavesurfer>

Multimodal speech synthesis - (NGSLT 2006 [83])

WaveSurfer Tool Demo



Multimodal speech synthesis - (NGSLT 2006 [84])

How to obtain data?

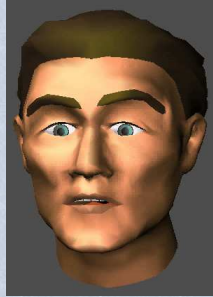


Qualisys recordings in Linköping


Multimodal speech synthesis - (NGSLT 2006 [85])

Combining model and data

Re-synthesis using speech movement recorded with Qualisys



Multimodal speech synthesis - (NGSLT 2006 [86])



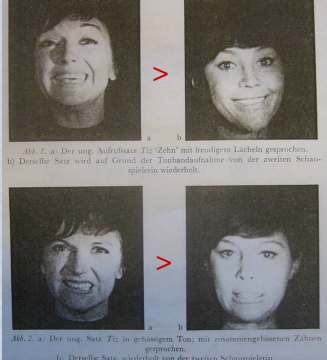
Preparing future multisensorial interaction research

1. technologies for speech-to-speech translation
2. detection and expressions of emotional states
3. core speech technologies for children

EU project: start October 2002, duration 2 YR
ITC-IRST (Trento) co-ordinates + 3*Germany
+ Italy + UK + Sweden
<http://pfstar.itc.it/>

Multimodal speech synthesis - (NGSLT 2006 [87])

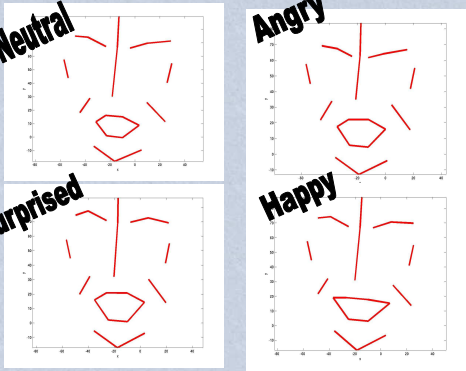
Vision from audio



original > mimic from audio

Fónagy, 1967 "Hörbare Mimik", *Phonetica*

Multimodal speech synthesis - (NGSLT 2006 [88])



Neutral

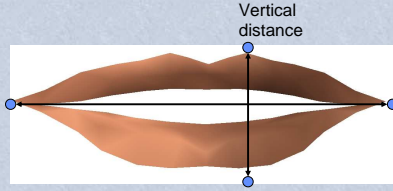
Angry

Surprised

Happy

Multimodal speech synthesis - (NGSLT 2006 [89])

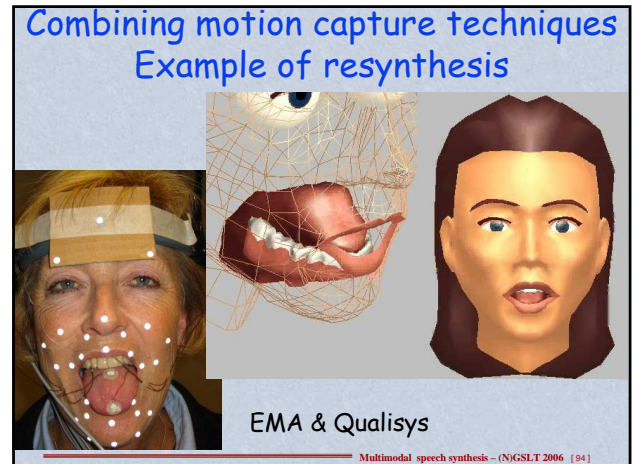
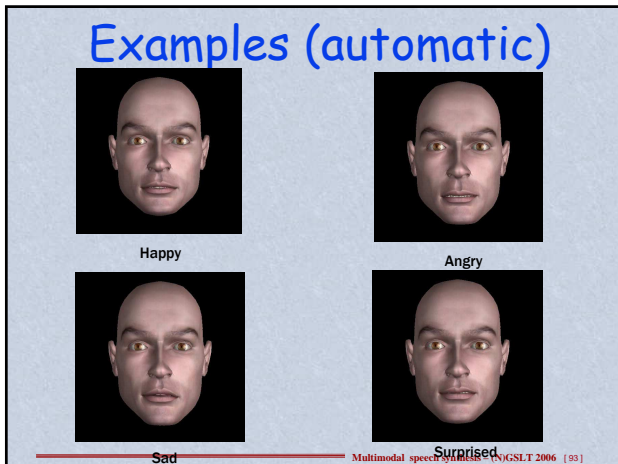
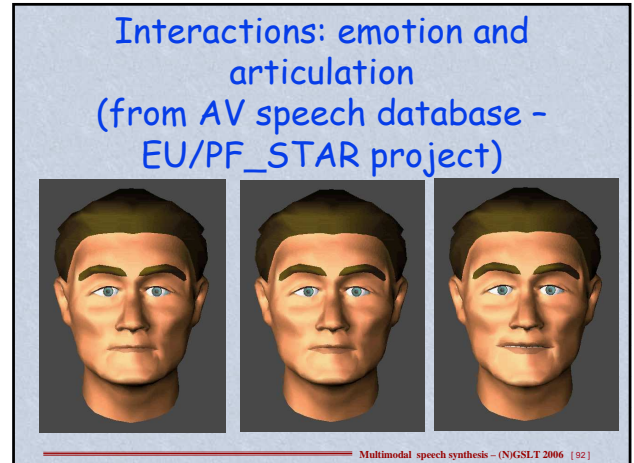
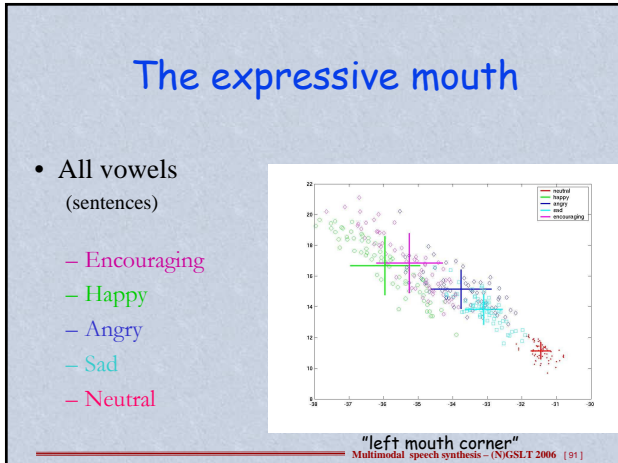
Measurement points for lip coarticulation analysis



Vertical distance

Lateral distance

Multimodal speech synthesis - (NGSLT 2006 [90])



Collection of audio-visual databases: interactive spontaneous dialogues

- Eliciting technique: information seeking scenario
- Focus on the speaker who has the role of information giver
- The speaker seats facing 4 infrared cameras, a digital video-camera, a microphone The other person is only video recorded.

Multimodal speech synthesis - (NGSLT 2006 [95])



Conversation with agent



Multimodal speech synthesis - (NGSLT 2006 [97])

Eyebrow vs intonation



- 1 No eyebrow motion
- 2 Eyebrow motion controlled by the fundamental frequency of the voice
- 3 Eyebrow motion at focal accents +
- 4 Eyebrow motion at the first focal accent +

“Jag heter Axel, inte Axell” (translation: “My name is Axel, not Axell”). In Sweden Axel is a first name as opposed to Axell, which is a family name.

Multimodal speech synthesis - (NGSLT 2006 [98])

Experiment

- Speech material
 - När pappa fiskar stör, piper Putte
When dad is fishing sturgeon, Putte is whimpering
 - När pappa fiskar, stör Piper Putte
When dad is fishing, Piper disturbs Putte
- 6 versions
 - 1 static, 5 eyebrow raising on successive content words
- 20 stimuli (6 x 3) plus first and last
- Subjects: 21 students at KTH
 - 14 native Swedish, 7 non-Swedish

Multimodal speech synthesis - (NGSLT 2006 [99])

Eyebrow movement

- Hand edited with a synthesis parameter editor
- 500 ms
 - 100 ms dynamic rise
 - 200 ms static raised
 - 200 ms dynamic lowering

Multimodal speech synthesis - (NGSLT 2006 [100])



Multimodal speech synthesis - (NGSLT 2006 [101])

TWO EXAMPLES

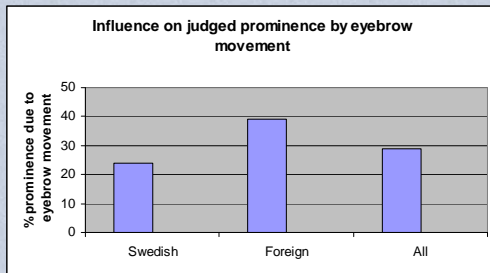


No eyebrow movement (neutral)

Eyebrow movement

Multimodal speech synthesis - (NGSLT 2006 [102])

Prominence increase due to eyebrow movement



Multimodal speech synthesis - (NGSLT 2006 [103])

Feedback experiment

- Mini dialogues (two turns)
- Travel agent application
- Both visual and acoustic feedback cues
- **Affirmative cues** – agent understands/accepts the request
- **Negative cues** – agent is unsure about the request (seeks confirmation)
- Six cues hypothesised

Granström, House & Swerts (2002)

Multimodal speech synthesis - (NGSLT 2006 [104])

Parameter settings to create different stimuli



	Affirmative setting	Negative setting
Smile	Head smiles	Head has neutral expression
Head movement	Head nods	Head leans back
Eyebrows	Eyebrows rise	Eyebrows frown
Eye closure	Eyes close a bit	Eyes open widely
F0 contour	Declarative intonation	Interrogative intonation
Delay	Immediate reply	Slow reply

Multimodal speech synthesis - (NGSLT 2006 [105])

Cue strength - demonstration

Human: "Jag vill åka från Stockholm till Linköping"

Agent: "Linköping"

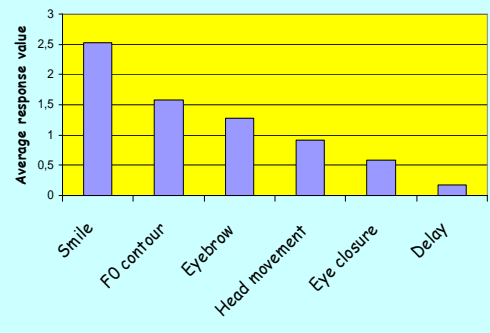
Affirmative cues:

- * None
- * All
- * Smile (only)
- * F0 contour (only)
- * No brow frown (only)
- * Head movement (nod) (only)
- * Eye closure (only)
- * Delay (only)
- * All



Multimodal speech synthesis - (NGSLT 2006 [106])

Cue strength



Multimodal speech synthesis - (NGSLT 2006 [107])

Conclusions

- Eyebrow movement can be an independent cue to prominence
- Non-native Swedish listeners rely more on the visual cues
- Interaction
 - visual and acoustic cues
 - visual cues and prominence expectation
- Further work on interaction
 - prominence, mood and attitude (demo)

Multimodal speech synthesis - (NGSLT 2006 [108])

Examples on the use of eyebrow and head motion (from the August dialogue system)



Translation: "Symmetrical works of art easily become dull just like symmetrical beauties; impeccable or flawless people are often unbearable." (Strindberg 1907)

Multimodal speech synthesis - (NGSLT 2006 [109])

Different characters



Multimodal speech synthesis - (NGSLT 2006 [110])

Talking heads on the Science Museum, Stockholm (Exhibition: FrittFram)



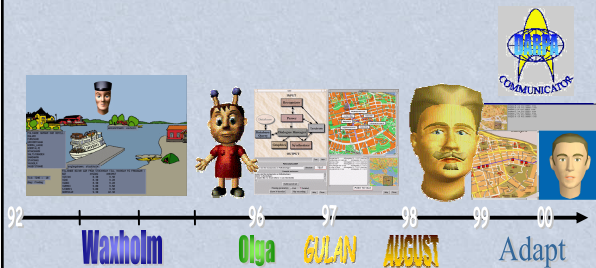
Multimodal speech synthesis - (NGSLT 2006 [111])

Speech synthesis applications

- Dialogue systems (at next NSGLT at KTH)
- Translators
- Aids for disabled
- Mediated communication
- New HMI
- Language learning
-

Multimodal speech synthesis - (NGSLT 2006 [112])

Dialog systems at KTH

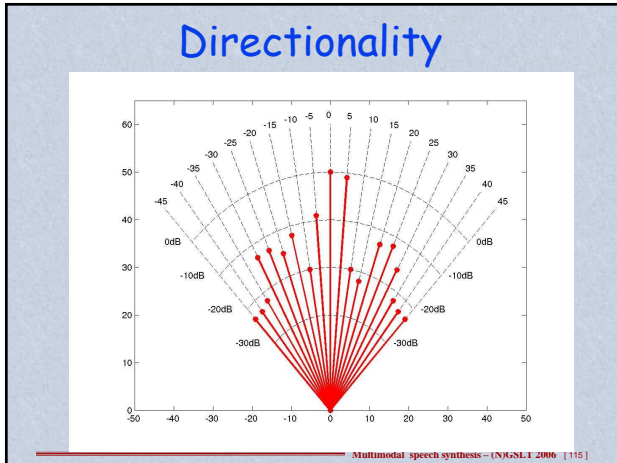


Multimodal speech synthesis - (NGSLT 2006 [113])

Targeted audio & talking head for personal announcements (EU/Chil project)



Multimodal speech synthesis - (NGSLT 2006 [114])



- ### Speech technology for visually impaired persons
- First synthesis application
 - Intelligibility vs. naturalness
 - Screen reader vs. GUI
 - Talking books and newspapers
 - “Design for all” or special demands
 - E.g. Rapid speech - 500 wpm
- Multimodal speech synthesis - (NGSLT 2006 [116])

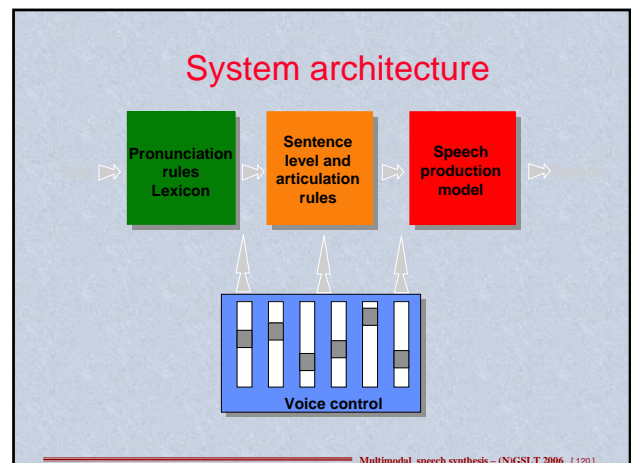
Aid for speech impaired The EU/Vaess project

Vaess - Voices, attitudes and emotions in speech synthesis
The Vaess Text-to-Speech converter enables speech impaired persons to communicate by individualised speech

Multimodal speech synthesis - (NGSLT 2006 [117])

- ### Synthesis Software Development
- Improved speech synthesis for disabled and elderly people
 - Danish, British English, Spanish and Swedish speech synthesis
 - Development of new voices
 - Development of different speaking styles
 - Experiments with new synthesis strategies for emotive speech
 - Uses an extension of the KTH/Infovox Speech Synthesizer
 - Based on analysis of human speech databases
- Multimodal speech synthesis - (NGSLT 2006 [118])

- ### User controlled “voice fitting”
- Direct access to selected voices
 - Individual settings easy to use
 - Phonetic rules use slide buttons as inputs
 - Synthesizer implemented with great flexibility
 - Examples of possible adjustments
 - Vocal tract size
 - Voice source characteristics
 - Pitch dynamics
 - Degree of clear or reduced speech
- Multimodal speech synthesis - (NGSLT 2006 [119])

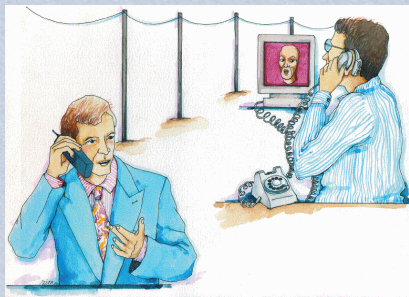


Speech synthesis for non-vocal persons
Cameleon CV - a speech prosthesis (EU/Vaess project)




Multimodal speech synthesis - (NGSLT 2006 [121])

Speech synthesis for hard of hearing persons
The Teleface application



Multimodal speech synthesis - (NGSLT 2006 [122])

"The TELEFACE project" (simulated)



Multi-modal Speech communication for the hearing impaired

Continues in EU project SYNFACE, aiming at a real-time demonstrator

Multimodal speech synthesis - (NGSLT 2006 [123])

EU-project Synface - Coordinated by KTH



<http://www.speech.kth.se/synface/>

Multimodal speech synthesis - (NGSLT 2006 [124])

Instruction video

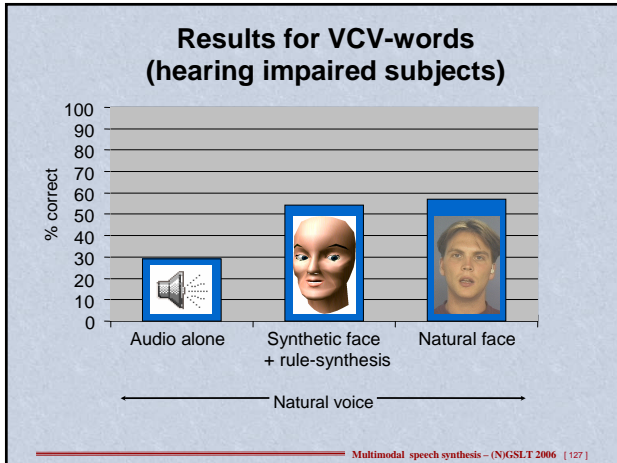


Multimodal speech synthesis - (NGSLT 2006 [125])

Formal intelligibility test

- Material: VCV (symmetric vowel context)
 - 2 vowels: /**ʊ**, **a**/
 - 17 consonants: /**p**, **b**, **m**, **f**, **v**, **t**, **d**, **n**, **s**, **l**, **r**, **k**, **ŋ**, **ʃ**, **ç**, **j**/
- Task: consonant identification
- Synthetic face with human speech
- hard of hearing subjects (or KTH students)
- Additive white noise, -3 dB SNR (if normal hearing)

Multimodal speech synthesis - (NGSLT 2006 [126])



Better than humans?

	bil	labd	den	pal	vel
bilabial	100				
labiodental	96,3	3,7			
dental	3,0	78,0	5,5	13,4	
palatal		9,9	70,4	19,8	
velar		4,9	16,0	79,0	

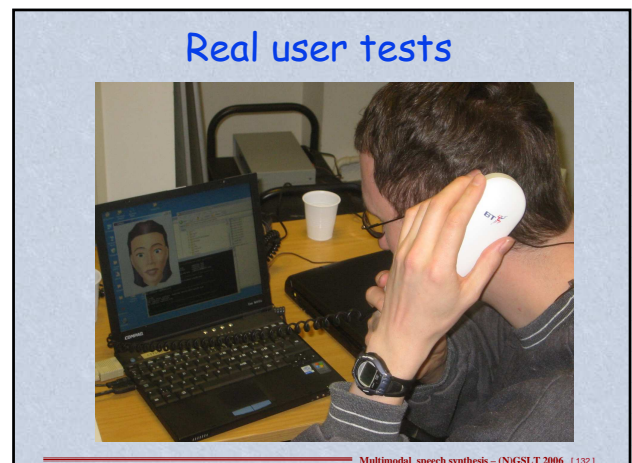
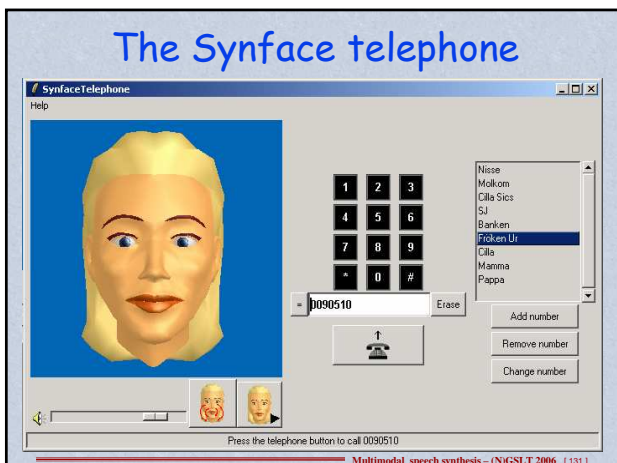
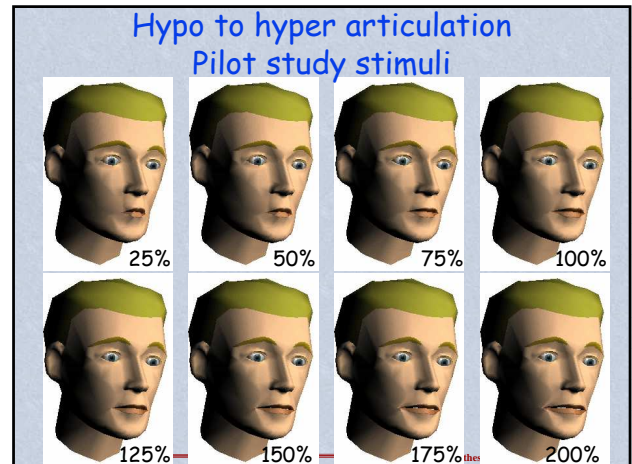
Synthetic face

	bil	labd	den	pal	vel
bilabial	96,3		2,5	1,3	
labiodental		92,6	5,6	1,9	
dental			85,8	7,4	6,8
palatal			1,2	17,3	71,6
velar				2,5	25,0

Natural face

Multimodal speech synthesis - (NGSLT 2006 [128])

- ### Possible improvements to "lip readability"
- Great variation in human speakers due to for example
 - Speaking rate
 - Extent of articulatory movements (the hypo - hyper dimension)
 - Anatomy, facial hair
 - Light, distance, viewing angle...
- Multimodal speech synthesis - (NGSLT 2006 [129])



Jonas Beskow - mottagare av Chester Carlsons forskningspris 2 februari, 2006



Multimodal speech synthesis - (NGSLT 2006 [133])

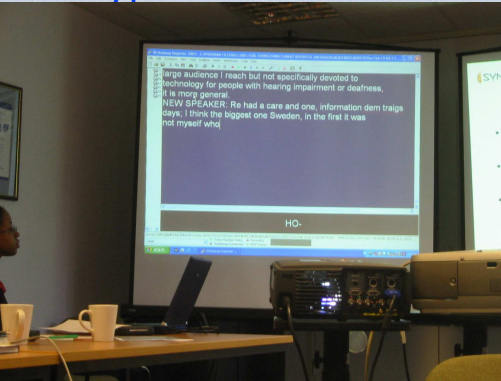
Talking faces for speech reading support in TV



2005 KTH Centrum för talteknologi

Multimodal speech synthesis - (NGSLT 2006 [134])

Texttype - man or machine?



Multimodal speech synthesis - (NGSLT 2006 [135])

Language learning

- Oral proficiency training
- Possible display of internal articulations
- Exploiting hyper/hypo dimension
- Training in dialogue context
- Always available conversational partner
- Untiring model of pronunciation
 - everything from phonemes to prosody

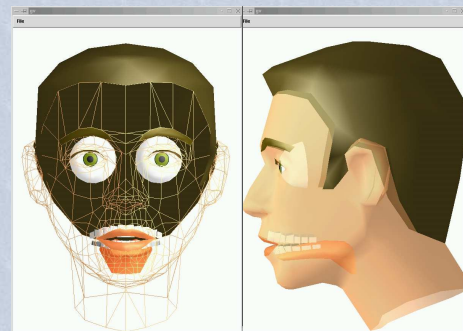
Multimodal speech synthesis - (NGSLT 2006 [136])

Automatic tutor simulation



Multimodal speech synthesis - (NGSLT 2006 [137])

Different representations



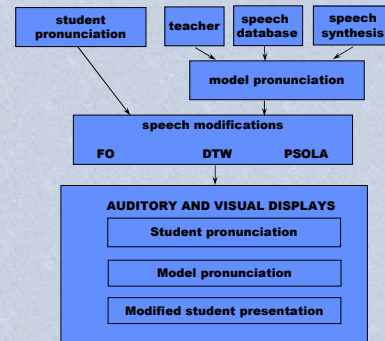
Multimodal speech synthesis - (NGSLT 2006 [138])

Speak&Spell - TexasInstrument Christmas 1978



Multimodal speech synthesis - (NGSLT 2006 [139])

Speech training using intonation models



Multimodal speech synthesis - (NGSLT 2006 [140])

Demo of prototype

Pohlman - weatherman from the south

“Sen drar hela det här moln- och regnområdet i alla fall vidare österut” (~then, this whole cloud and rain system moves eastward)

1 Original recording

“Teacher”(sound only)-original-modified

2 Stockholm

3 South Swedish

4 Synthesis

Multimodal speech synthesis - (NGSLT 2006 [141])

1 Original recording

2 Stockholm “Teacher”(sound only)-original-modified

3 South Swedish

4 Synthesis



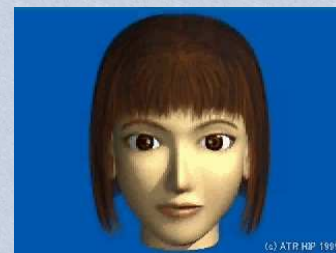
Multimodal speech synthesis - (NGSLT 2006 [142])

Articulatory training

- Stylized
- Program Fonem -Johan Liljencrants

Multimodal speech synthesis - (NGSLT 2006 [143])

Reiko Yamada ATR, 1999



(c) ATR 1999


Multimodal speech synthesis - (NGSLT 2006 [144])

new national project ARTUR

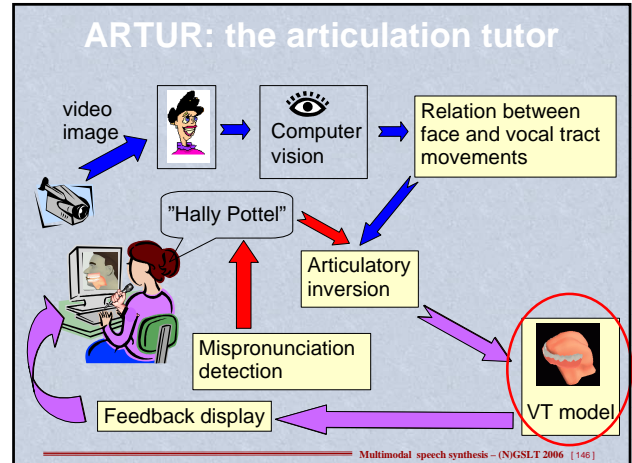
What?
Automatic articulatory feedback display using face and vocal tract models.

For whom?
Hearing impaired children, second-language learners, speech therapy patients.

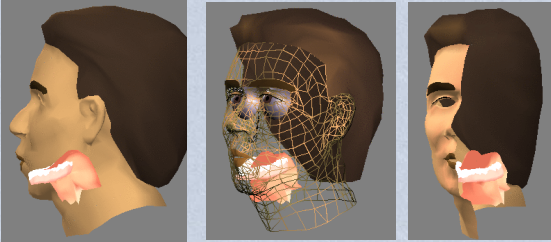
How?
Contrasting the user's articulation with a correct one.



Multimodal speech synthesis - (NGSLT 2006 [145])



3D tongue movements




[s] [k] vowels

“OK” 3D tongue movements (difficult to define an objective evaluation criteria - suggestions?)
The tongue movements can be presented from different views

Multimodal speech synthesis - (NGSLT 2006 [147])

CTT Virtual Language Tutor

- Practice dialogues
- Correct your pronunciation
- Keep track of your improvements
- Tailor lessons based on your interaction




Multimodal speech synthesis - (NGSLT 2006 [148])

CTT Virtual Language Tutor

Different Types of Users:

- Swedish children learning English
- Adult immigrants learning Swedish
- Adult Swedes wanting to improve aspects of English (e.g. corporate English, technical English)
- Native Swedes with language disabilities wanting to improve their Swedish



Multimodal speech synthesis - (NGSLT 2006 [149])

Discriminate acceptable from unwanted variation

- How to do it (automatically)
- What are the aims of L2 learning
 - Less accentedness
 - Comprehensibility
 - Intelligibility
 - More? – Acceptability in context
- Economy of language learning
- Could a virtual tutor help?

Multimodal speech synthesis - (NGSLT 2006 [150])

The cost of pronunciation errors

- Video

Multimodal speech synthesis - (NGSLT 2006 [151])

The virtual teacher vision already at STiLL, Marholmen 1998!

ISCA STiLL/InStil sig revival?

Speech Technology
in Language Learning



Unacceptable pronunciation needs to be identified

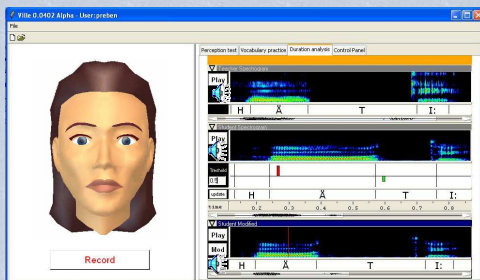


from a Vinnova video (on CTT's webpage) (NGSLT 2006 [153])

First tried at the VISPP summerschool in Palmse, August 2005 (for Estonian)



Nordic project -NordPlus Sprog Using VILLE - CTT Virtual Language Tutor



Multimodal speech synthesis - (NGSLT 2006 [155])



HAL'S LEGACY: 2001'S COMPUTER AS DREAM AND REALITY

Chapter 6
"The Talking Computer": Text to Speech Synthesis
Joseph P. Olive
Chapter 7
When Will HAL Understand What We Are Saying?
Computer Speech Recognition and Understanding
Raymond Kurzweil

<http://mitpress.mit.edu/e-books/Hal/>

Multimodal speech synthesis - (NGSLT 2006 [156])

