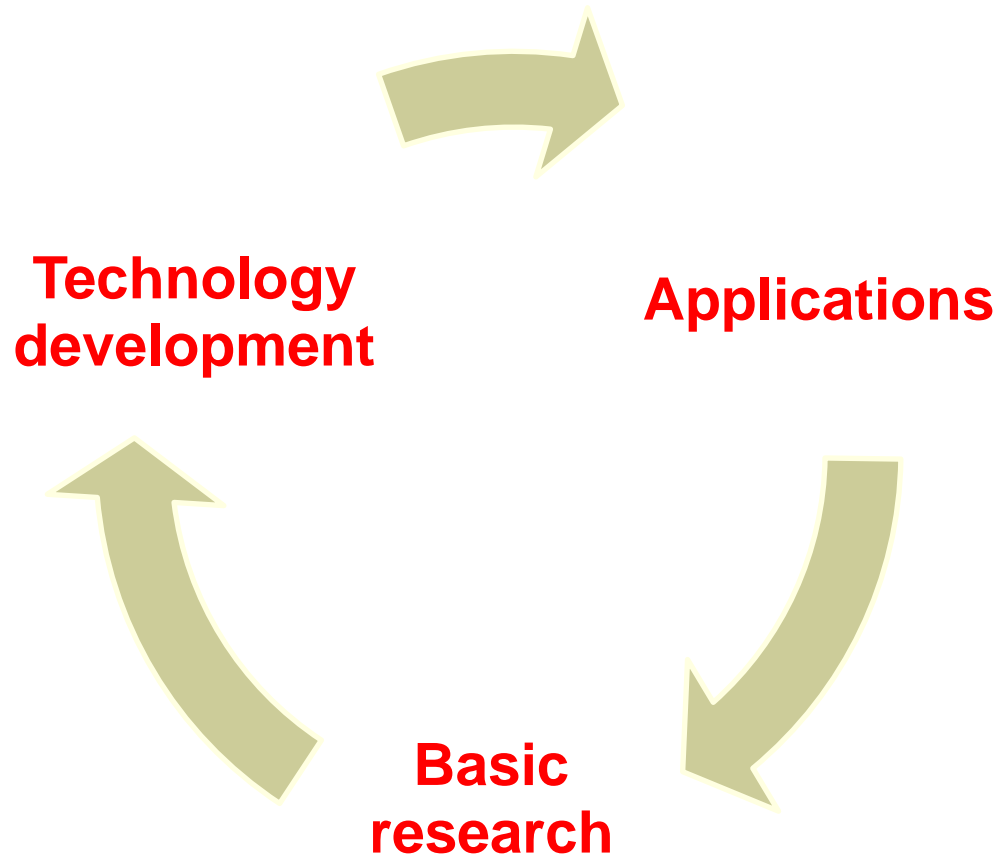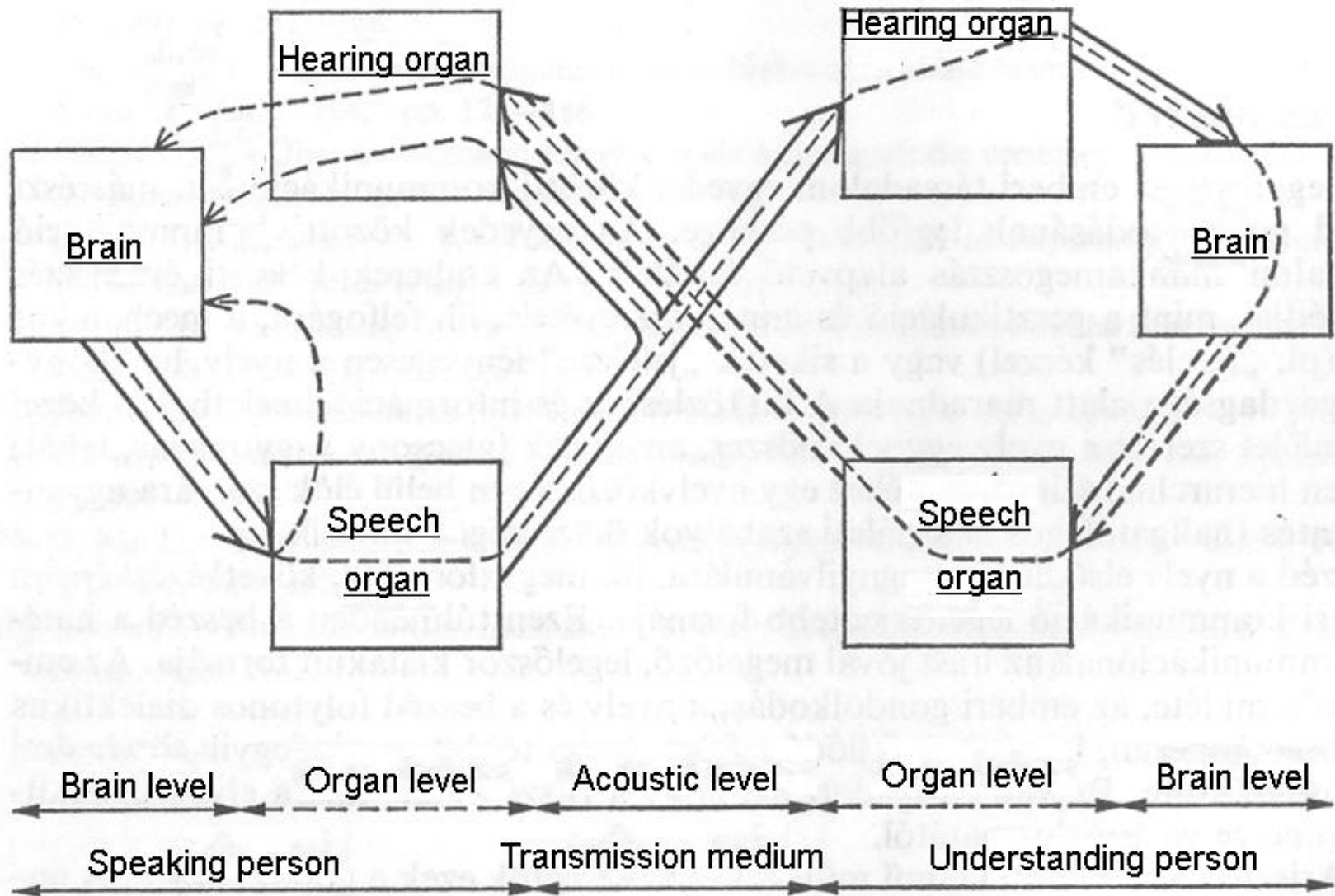# 10. Speech response systems

# Topic

- The creation of machine generated speech modern / contemporary technologies
- Modeling approaches
    - Generative models (articulation, source-filter, …)
    - Production models (concatenation, …)
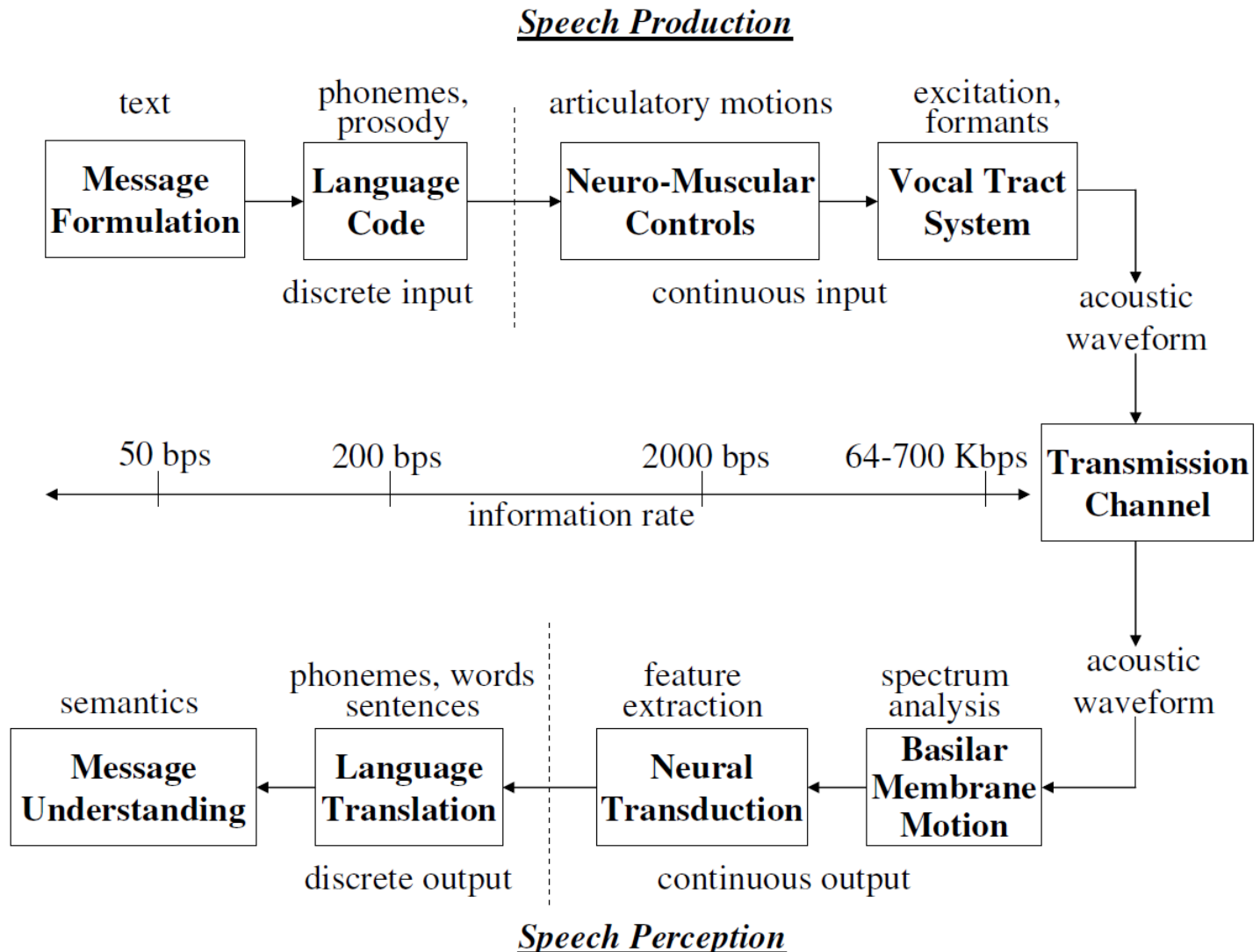- Technological solutions
- Applications
- Outlook

SmartLab
Intelligent Interactions
smartlab.tmit.bme.hu

BME TMIT

NVIDIA
GPU
EDUCATION
CENTER

2/62

# The research cycle of speech technology



**Applications**

**Technology development**

**Basic research**

# The natural speech chain

# The natural speech chain



*Speech Production*

text — **Message Formulation** — phonemes, prosody — **Language Code** — articulatory motions — **Neuro-Muscular Controls** — excitation, formants — **Vocal Tract System**

discrete input — continuous input — acoustic waveform

information rate: 50 bps | 200 bps | 2000 bps | 64-700 Kbps — **Transmission Channel**

acoustic waveform

semantics — **Message Understanding** — phonemes, words sentences — **Language Translation** — feature extraction — **Neural Transduction** — spectrum analysis — **Basilar Membrane Motion**

discrete output — continuous output

*Speech Perception*

**SmartLab**
Intelligent Interactions
smartlab.tmit.bme.hu

BME TMIT

NVIDIA.
GPU
EDUCATION
CENTER

5/62

# The natural speech chain

SmartLab
Intelligent Interactions
smartlab.tmit.bme.hu

BME TMIT
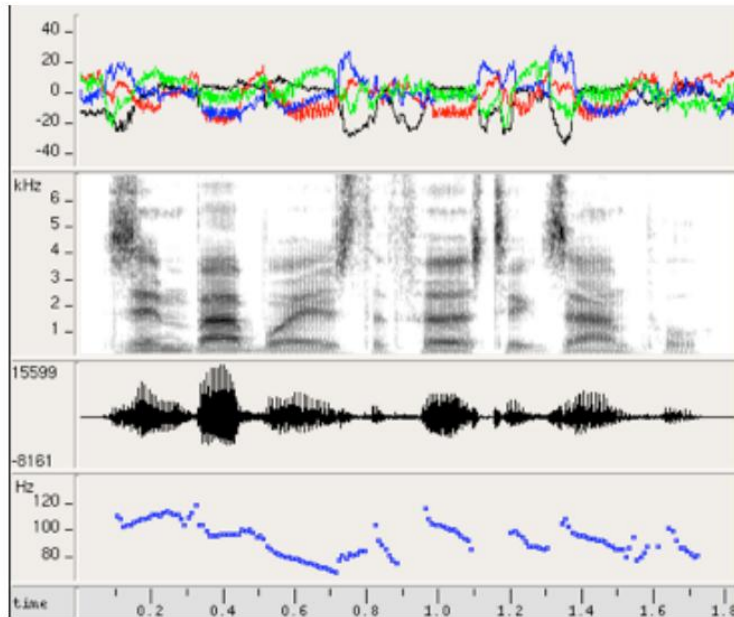
NVIDIA.
GPU
EDUCATION
CENTER

6/62

# Machine speech generation

- Machine speech generation ( speech synthesis ???)
  - Playback or concatenation based on simple symbols
    - Gramophone
    - Tape recorder
    - CDs …
  - From text ( speech synthesis ???)
    - Text-to-speech (TTS)
    - Text-to-speech with additional information ( concept-to-speech - CTS)
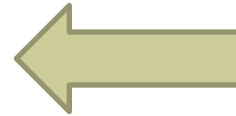  - Audiovisual speech synthesis

**SmartLab**
Intelligent Interactions
smartlab.tmit.bme.hu

BME TMIT

NVIDIA.
GPU
EDUCATION
CENTER

# Modeling approach
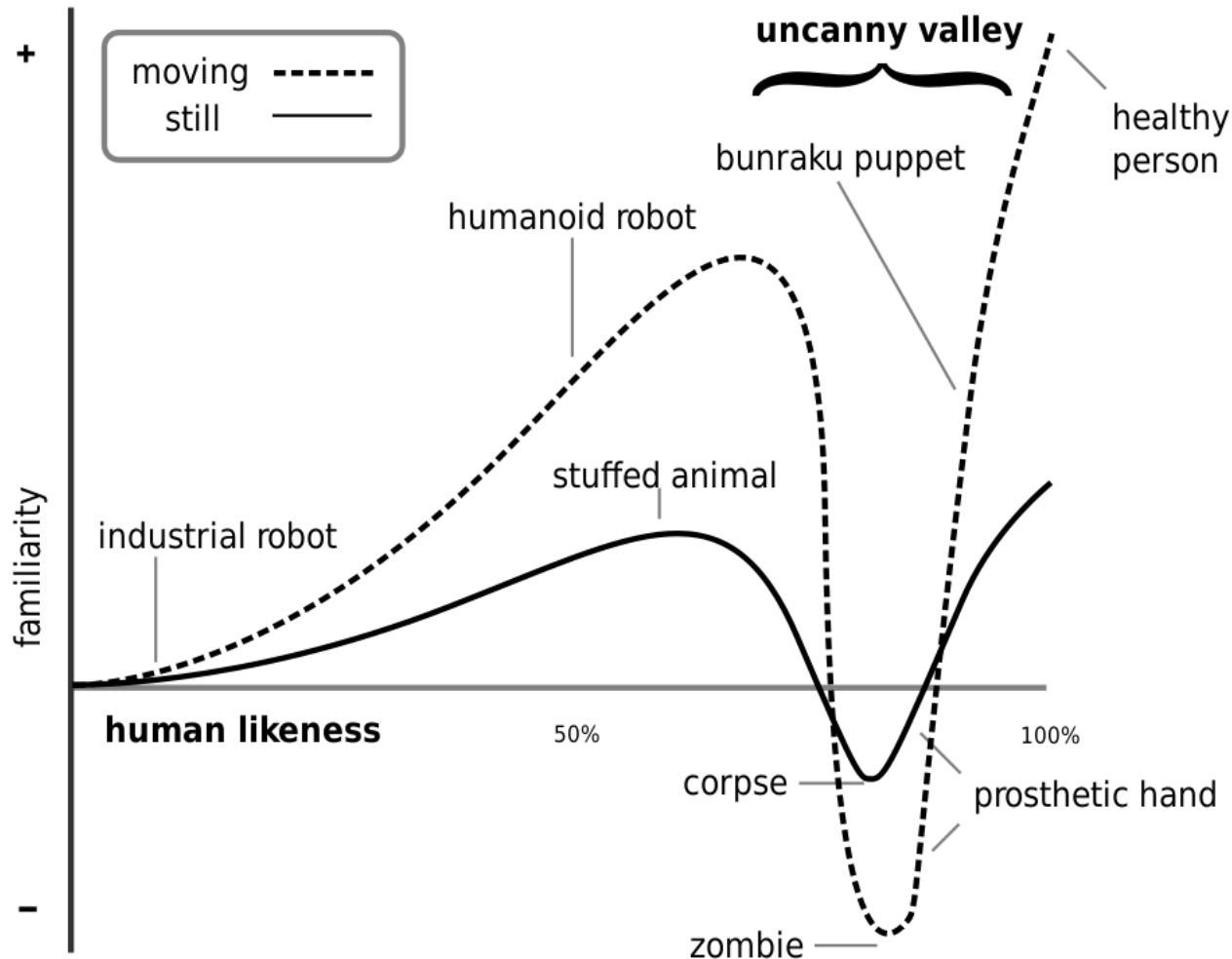


Speech in the acoustic domain

*Production model*
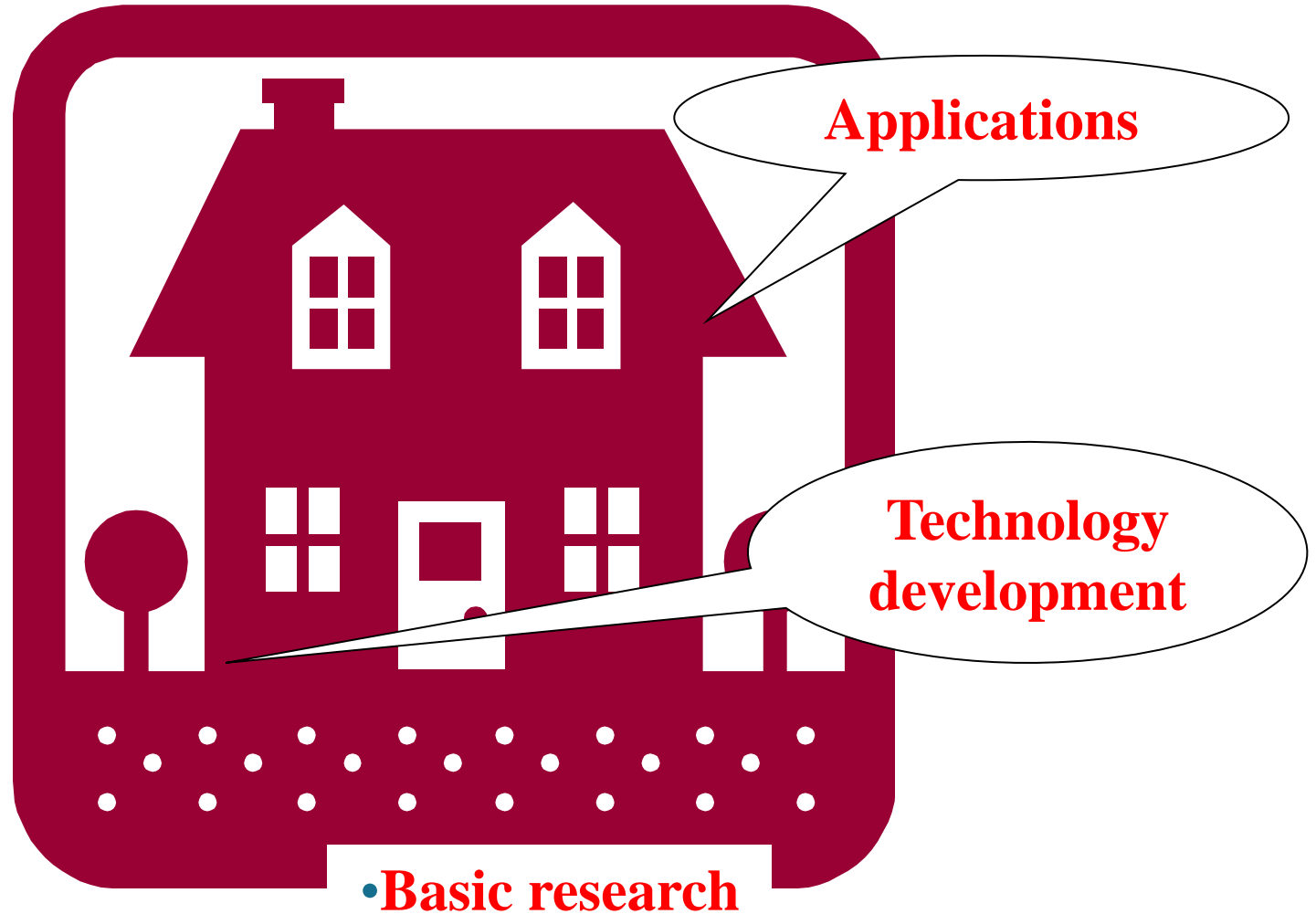
Speech in the articulatory range

*Generative model*

# Research challenge

## Avoiding the uncanny valley

# Speech technology research cycle



**Applications**

**Technology development**

•**Basic research**

# Modeling approach

- Generative, articulation model
  - Segmental level
    - intelligibility of sounds
  - Suprasegmental level
    - correct intonation (prosody)
  - Non-linguistic features
    - individual tone, emotion, …

*SmartLab*
Intelligent Interactions
smartlab.tmit.bme.hu

BME TMIT

NVIDIA.
GPU
EDUCATION
CENTER

11/62

# Kempelen's talking machine (1791)



**RECONSTRUCTED, WORKING VERSION (2001 MTA NyI , 2020 BME)**

# Kempelen's talking machine (1791)

# Kempelen's talking machine (1791)

- **Resources**
  - Mechanical articulation model
  - Direct control of model elements by human hands
- **Solution**
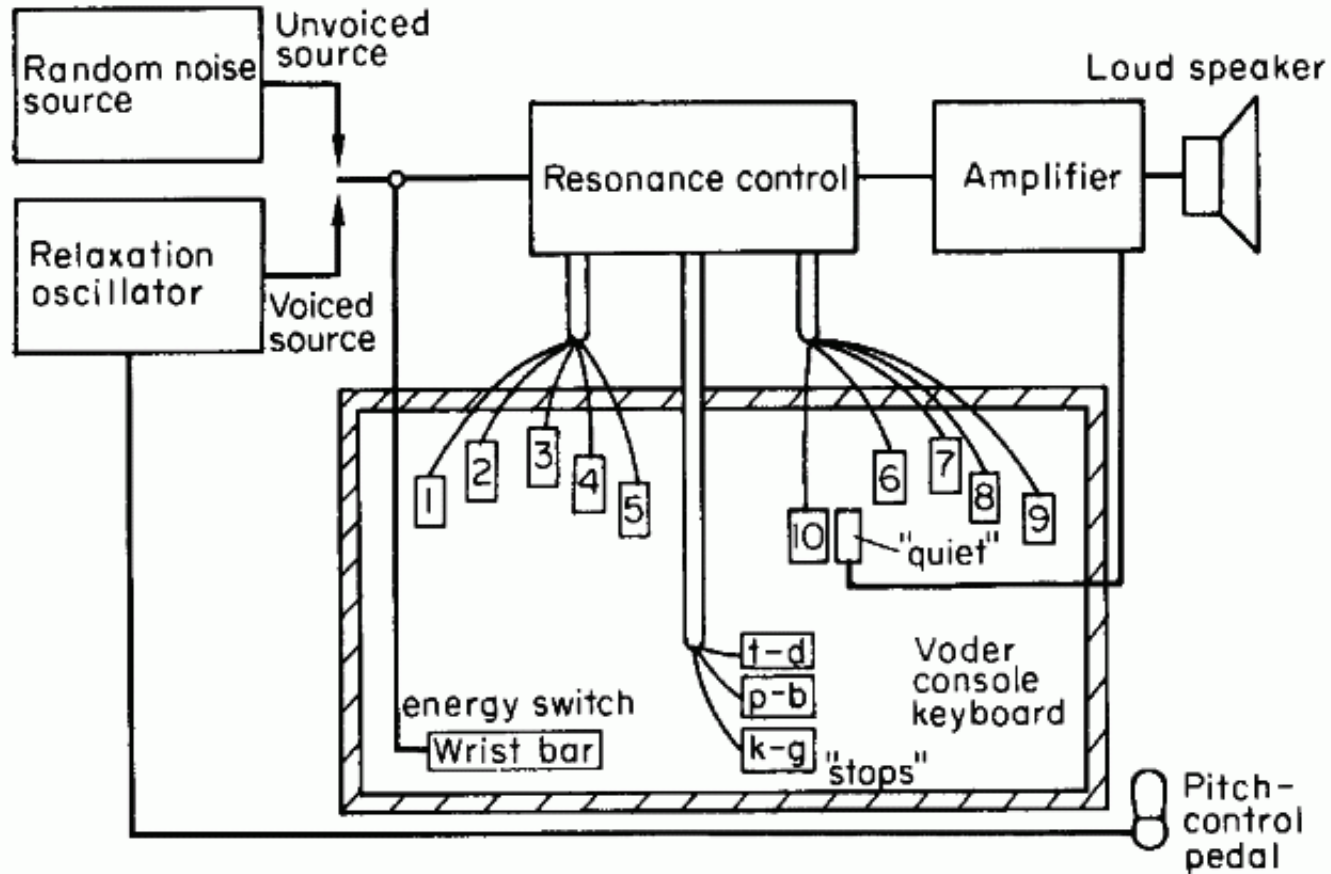  - Continuous / sustained sounds
  - "simple" sound concatanations

*SmartLab*
Intelligent Interactions
smartlab.tmit.bme.hu

BME TMIT

NVIDIA.
GPU
EDUCATION
CENTER

14/62

# Kempelen's "choir" of machines



**Source: Meeting of Kempelen machines, September 2019, Vienna**

# Dudley-Voder (1939)

**SmartLab**
Intelligent Interactions
smartlab.tmit.bme.hu

BME TMIT

NVIDIA.
GPU
EDUCATION
CENTER

16/62

# Dudley-Voder (1939)

**SmartLab**
Intelligent Interactions
smartlab.tmit.bme.hu

BME TMIT

NVIDIA.
GPU
EDUCATION
CENTER

17/62

# Dudley-Voder (1939)

■ Resources

   ■ Electro -mechanical articulation model

   ■ Keyboard control with human hand

■ Solution

     ■ In principle, a loose vocabulary

     ■ Demonstration trick

*SmartLab*
Intelligent Interactions
smartlab.tmit.bme.hu

BME TMIT

NVIDIA.
GPU
EDUCATION
CENTER

18/62

# Serial and parallel formant synthesizer (1983-84)

# Holmes Olaszy Klatt
## (1961 -1982 -84)

- Resources
  - Electronic formant/articulation model
    (equivalent LPC representation solved)
  - Rule-based processing
  - Control with computer
- Solution
  - In principle, a loose vocabulary
  - Singing
  - Small acoustic database (up to 1kByte)

# Electronic formant synthesis solutions

## KlatTalk/DecTalk 1982





## HungaroVox 1982

## MultiVox 1986-2002

# Modeling approach

- Production model
  - Audio elements of varied details
    - Concatenation (modified by signal processing)
  - Segmental level
    - intelligibility of sounds
  - Suprasegmental level
    - correct intonation (prosody)
  - Non-linguistic features
    - individual tone, emotion, …

**SmartLab**
Intelligent Interactions
smartlab.tmit.bme.hu

BME TMIT

NVIDIA
GPU
EDUCATION
CENTER

22/62

# Miklós Bánó's talking machine (1916)

- June 21 , 1919. Hungarian Patent Office Office , Patent description , number 74361 , class IX/d

  *A talking machine suitable for reading any text*

- Dr. Miklós Bánó certified engineer and economic engineer In Budapest
- The submission day:

  ## November 30, 1916

*SmartLab*
Intelligent Interactions
smartlab.tmit.bme.hu

BME TMIT

NVIDIA.
GPU
EDUCATION
CENTER

23/62

# Content of the patent
## The Machine Reader (1916)

■ 5 claims

  ■ 1. "A speaking machine capable of reproducing any text, characterized by recording individual sounds or groups of sounds separately and by playing these separately recorded sounds or groups of sounds in the order corresponding to them."

  ■ 2. "An embodiment of the machine protected against the claim of 1., characterized in that each sound of the language is recorded on a separate recording cylinder or disc"

# Content of the patent
## The Repentant (1916)

- 3. sound storage in a cylinder .

- 4. 1. + 3.

- 5. " Under Claim 1 protected machine implementation shape , characterized by the same recording cylinder or discs activation by typewriter-style buttons which can generate speech."

- ( Naive ) Example : per cylinders of *r+o+z* activated after each other melt into the word *rose*

# Application example
## Bánó (1916)

- "Connected to a typewriter...suitable for the visually impaired typist to listen to the written text"

- "The speech impaired typist ... should be able to communicate what he has to say aloud with the help of the device"

- ..."also suitable for the reproduction of prescribed or duplicated texts"

# Miklós Bánó's talking machine (1916)

- Resources
  - Electro-mechanical speech concatenation model
  - Keyboard control with human hand
- Solution
  - No working sample
    - Sustained speech sounds
    - In principle, unlimited vocabulary
    - Simulation reconstruction (2010)

# Overview of machine speech generation methods

| Speech synthesis method | Prosody production | Speech database type |
|---|---|---|
| "classical" formant synthesis (the starting method) | rule-based, with encoder control parameters | parametric (formant filter model) |
| element-concatenation (diad/diphone) | rule-based, with waveform modification | sound, diad/diphone waveform elements (logatoms) |
| element-concatenation ( triad/triphone ) | rule-based, with waveform modification | sound, triad/triphone and diad/diphone waveform elements (logatoms) |
| element selection (corpus) | indirect, pattern search based on the time scale of the current sentence, typically without waveform modification | large waveform database (from readings) of variable-sized elements (words, word strings, sentences, etc.) |
| statistical parametric | with a statistical (HMM or DNN) model, which operates based on parameter n-grams at the sentence level | parametric (LPC, harmonic+noise, sinusoidal, etc.) |
| waveform-based statistical parametric ( WaveNet /DNN) | with statistical (DNN) model | neural network parameters (learning from waveform and direct generation) |

# Database features

| Sign | Audio recording length/database length (minutes) | Gender | Language | Goals |
|------|--------------------------------------------------|--------|----------|-------|
| DIAD1- DIAD4 | 28/2.5 | man | Hungarian | Element-concatenation (diad-triad / diphone-triphone) research |
| DIAD5–DIAD8 | approx. 28 minutes/2.5 minutes | woman | | |
| TRIAD1 | approx. 120min/32 min | man | | |
| TRIAD2 | approx. 120min/32 min | woman | | |
| TIME1 | 630 minutes | woman | Hungarian | Item selection (corpus-based) research |
| FON1 | 100 minutes | | | |
| PAGE1 | 110 minutes | | | Experimental station passenger information system |
| ARU1 | 330 minutes | | | |
| UGYF1 | 505 minutes | | | |
| SZAM1 | 10 minutes/XXXX | | | Reading numbers up to 1 billion |
| SZAM2 | 10 minutes/XXXXX | man | | |
| RADIO | 516 minutes | 3 men | | |
| FON2-5 | approx. 130 minutes/person | 4 women | | Statistician |
| FON6-10 | approx. 130 minutes/person | 5 men | | parametric (HMM and |
| BEA1 | 30 minutes | woman | | DNN) research |
| BEA2 | 31 minutes | man | | |
| GABOR | 3 minutes | man | | Spemoticon research |

# Element concatenation model

Transcription rules,
exception dictionary

Waveform element
database

Physical rules

| phonetic transcription | → | prosody regulation | → | element-concatenation | → | prosody modification - signal processing | → |

input text

text,
phoneme
sequence

phoneme series,
I, $F_0$ , durations

segmental signal,
$F_0$ , durations

speech
signal

# Diad- triad concatenation ProfiVox (1999-)



papa = _p  pa  ap  pa a_

or

_p  pap  pa_

diad + triad

# Hungarian sound definition for synthesis

- Hungarian speech with 39 speech sounds
  - 24 consonants (C)
  - 14 vowels (V)
  - plus the pause ( _ sign) can be covered
    - long consonants are generated from their short versions using signal processing.
    - We are very sensitive to the quality of vowels, so it is advisable to treat short and long versions separately.

# Speech sound set

| No. | Hungarian character | Example | Own character set | IPA Unicode | IPAASCII |
|---|---|---|---|---|---|
| 1 | (pause) | | _ | | |
| 2 | á | láb | A: | 0250 02D0 | a&: |
| 3 | a | hat | A | 0254 | c& |
| 4 | o | sok | O | 006F | O |
| 5 | u | fut | U | 0075 | U |
| 6 | ü | süt | U | 0079 | Y |
| 7 | i | hit | i | 0069 | I |
| 8 | é | méz | E: | 0065 02D0 | e: |
| 9 | ö | köt | O | 00F8 | o/) |
| 10 | e | vet | e | 025B | E |
| 11 | b | bál | b | 0062 | B |
| 12 | p | tár | p | 0070 | P |
| 13 | d | dán | d | 0064 | D |
| 14 | t | tár | t | 0074 | T |
| 15 | g | gát | g | 0067 | G |
| 16 | k | kád | k | 006B | K |
| 17 | gy | gyár | G | 025F | j- |
| 18 | ty | tyúk | T | 0063 | C |
| 19 | m | már | m | 006D | M |
| 20 | n | nád | n | 006E | N |
| 21 | ny | nyom | N | 0272 | nj) |
| 22 | j | Jön | j | 006A | J |
| 23 | H | Hát | h | 0068 | H |
| 24 | v | Vád | v | 0076 | V |
| 25 | f | Fát | f | 0066 | F |
| 26 | z | Zár | z | 007A | Z |
| 27 | sz | Szép | s | 0073 | S |
| 28 | c | Cél | c | 0074 0073 | Ts |
| 29 | zs | Zsír | Z | 0292 | 3" |
| 30 | s | Só | S | 0283 | S |
| 31 | cs | Cső | C | 0074 0283 | TS |
| 32 | l | Láp | l | 006C | L |
| 33 | r | Rák | r | 0072 | R |
| 34 | ó | Pók | o: | 006F 02D0 | o: |
| 35 | ú | Kút | u: | 0075 02D0 | u: |
| 36 | ű | Fűt | U: | 0079 02D0 | y: |
| 37 | í | Szít | i: | 0069 02D0 | i: |
| 38 | ő | Sőt | O: | 00F8 02D0 | o/): |
| 39 | j* | Kapj | j | 006A | J |

# The 5 diads of the word alma



- The boundaries of the sound periods are shown in blue,

- the sound boundary in red

- the boundaries of the diads are the lines next to the markings in the lower gray bar

# The number of diad variants required for the synthesis of the Hungarian language

| Audio | CV | VC | CC | VV | _V and V_ | _C and C_ | Total |
|---|---|---|---|---|---|---|---|
| Quantity | 336 | 336 | 576 | 196 | 28 | 48 | 1520 |

# Prosodic units

# F0/sound duration modification
# PSOLA
## Pitch Synchronous OverLap-Add

# F0/sound duration modification
# PSOLA



https://wiki.aalto.fi/pages/viewpage.action?pageId=155477136

# Diad- triad concatenation ProfiVox (1999-)

- **Resources**
  - Digital diad - triad voice concatenation model
  - Rule-based processing
  - Control with computer

- **Solution**
  - Practically unlimited vocabulary
  - Flexible prosodic control
  - Acoustic database 1.5-60Mbyte

**SmartLab**
Intelligent Interactions
smartlab.tmit.bme.hu

BME TMIT

NVIDIA
GPU
EDUCATION
CENTER

# Stephen Hawking's machine voice (English - Hungarian)

## Dectalk 1982





## ProfiVox 2000 – 2014

Hungarian dubbing: Theory of Everything movie (Theory of Everything)

**SmartLab**
Intelligent Interactions
smartlab.tmit.bme.hu

BME TMIT

NVIDIA
GPU
EDUCATION
CENTER

# Corpus approach ( 2003- )

- N * Gbytes of storage space available
  - Store as many items as possible
  - Limit: limits of human reading with the same tone

- Corpus-based concatenation synthesis
  ( corpus-based , unit-selection synthesis)
  - Prosody can be controlled better indirectly
  - Exceptions are difficult to implement

# Corpus-based, waveform element selection speech synthesizer model

transcription rules,
exception dictionary

situation, rules
and/or speech
corpus

speech corpus

| phonetic transcription | → | prosody regulation | → | element selection and concatenation | → | (sometimes prosody modification ) | → |

input text

text,
phoneme
sequence

phoneme series,
symbolic or
physical prosody

concatenated
speech signal
(physical prosody)

speech
signal

# Indirect prosodic model

# Corpus approach ( 2003- )



Corpus-based concatenation

# Corpus approach
# ( 2003- )

- Resources
  - Digital corpus voice concatenation model
  - Statistical processing
  - Control with computer
- Solution
  - Optimally limited to a topic
  - Practically unlimited vocabulary with slightly worse sound quality
  - Acoustic database n*10Mbyte – n*Gbyte

**SmartLab**
Intelligent Interactions
smartlab.tmit.bme.hu

BME TMIT

NVIDIA
GPU
EDUCATION
CENTER

46/47

# A new voice at Hungarian train stations



MAGYARORSZÁG SZEMÉLYSZÁLLÍTÁSI VASÚTI TÉRKÉPE - 2011

Készítette: Balla László 2011.II. - www.vasut.info; www.vasut.tk

http://smartlab.tmit.bme.hu

# Hidden Markov Model (HMM)

# HMM evolution

**Speech recognition technology development (1970-90)**

**Dictation applications (1990-)**
**Dragon Dictate 1997,**
**Excessive Marketing promises 2000-**

**Markov**
**(end of the 19th century)**
**Signal processing models,**
**Algorithms 1970-80ies**

**Further development , other areas (1990-)**

# HMM features
## Tokuda et al . (1996-), Tóth, Németh (2008-)

■ Resources

- Digital parametric speech coding model
- Statistical processing (training database)
- Computer controlled process

■ Advantages

- Practically unlimited vocabulary
- Speaker adaptation
- Slightly poorer speech quality (coding distortions)
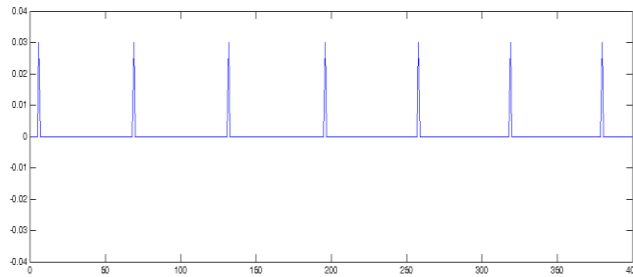- n*Mbyte acoustic database

**SmartLab**
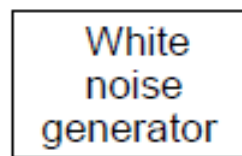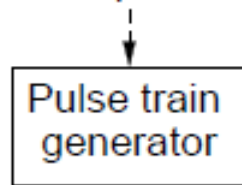Intelligent Interactions
smartlab.tmit.bme.hu

BME TMIT

NVIDIA.
GPU
EDUCATION
CENTER

# HMM-based TTS

# Spectrum modeling

SmartLab
Intelligent Interactions
smartlab.tmit.bme.hu

BME TMIT

NVIDIA.
GPU
EDUCATION
CENTER

52/47

# Speech signal excitation models
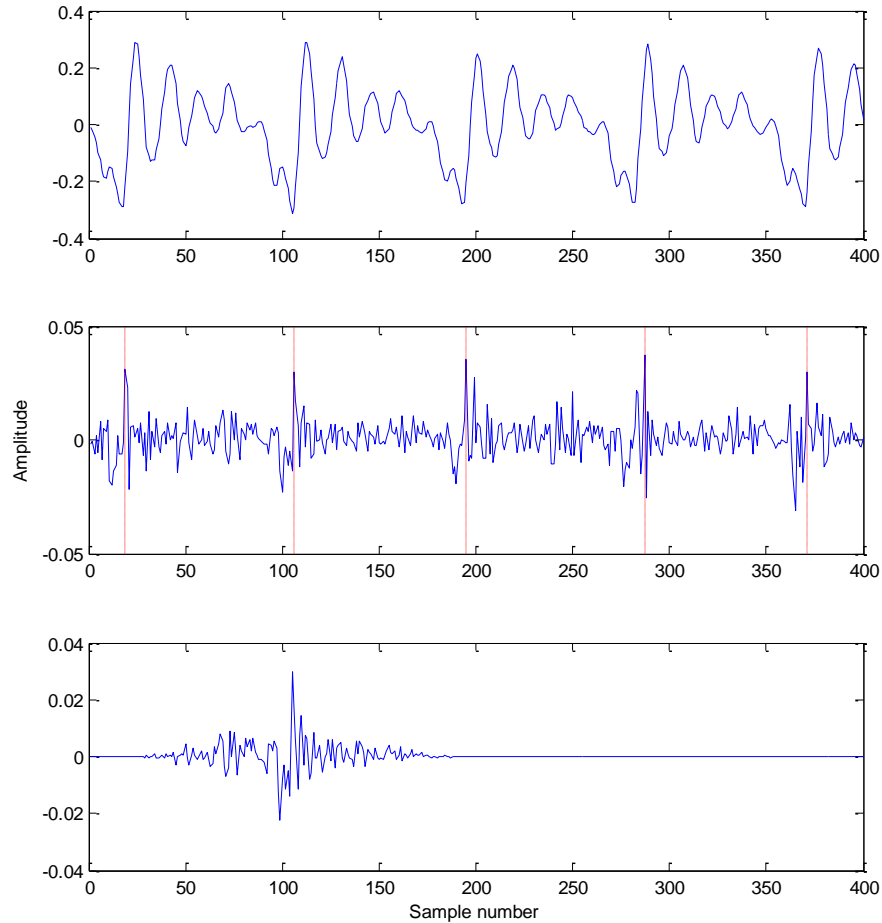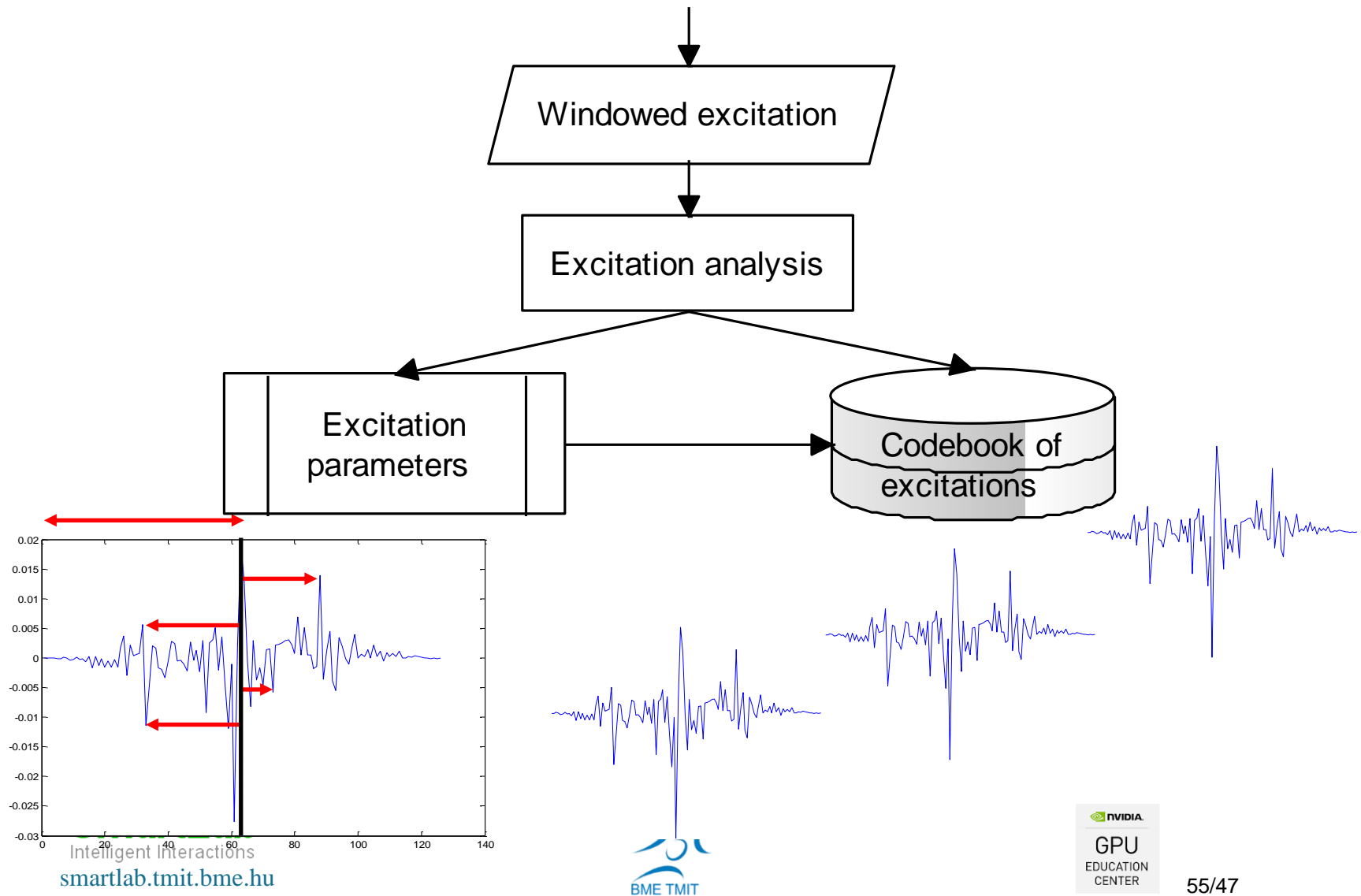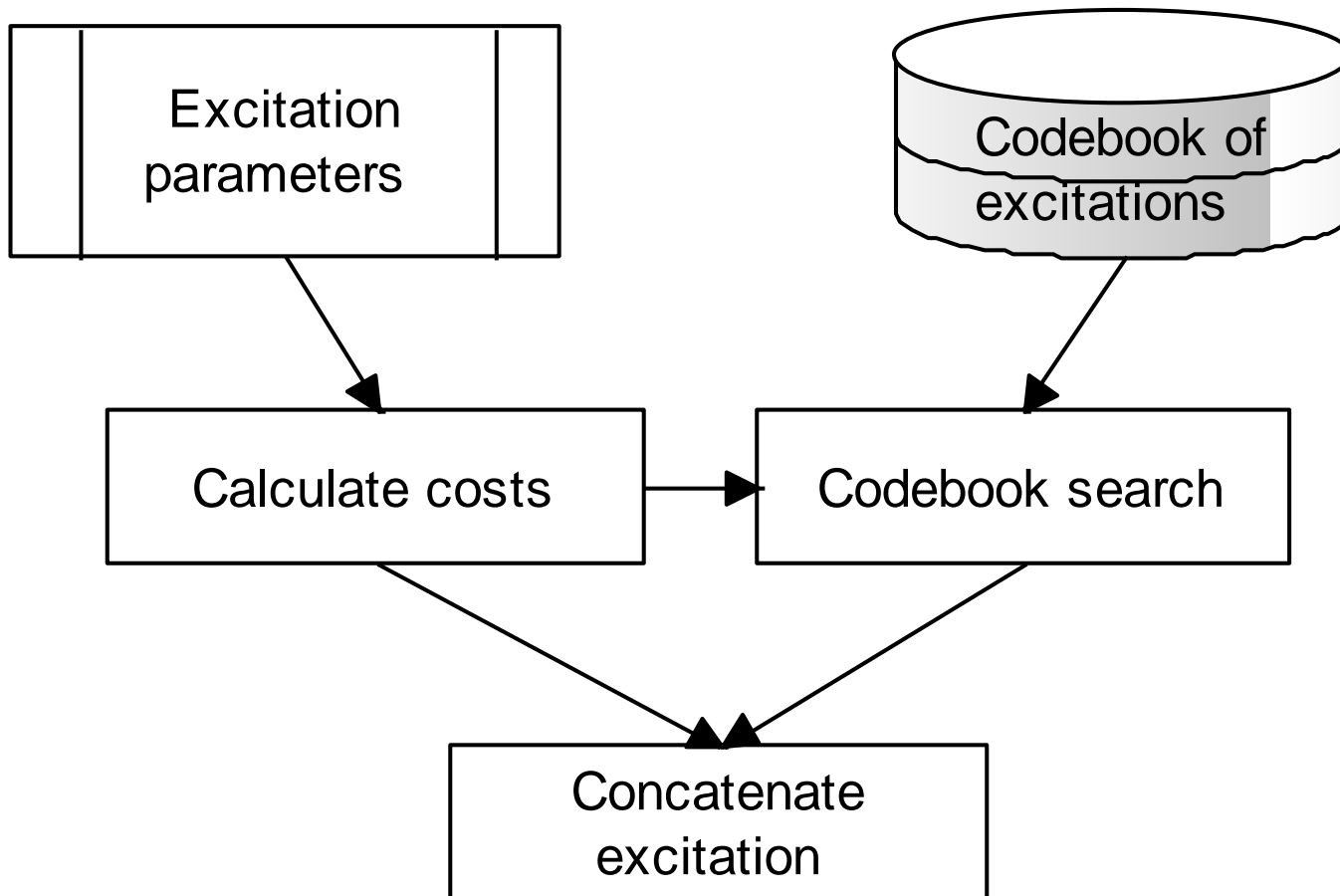
# Excitation signal from natural speech

# Excitation codebook ( codebook )

# Speech production

# Articulatory/formant rebirth (2009-)

- **Resources**
  - Hybrid approach
  - Digital articulation/formant model
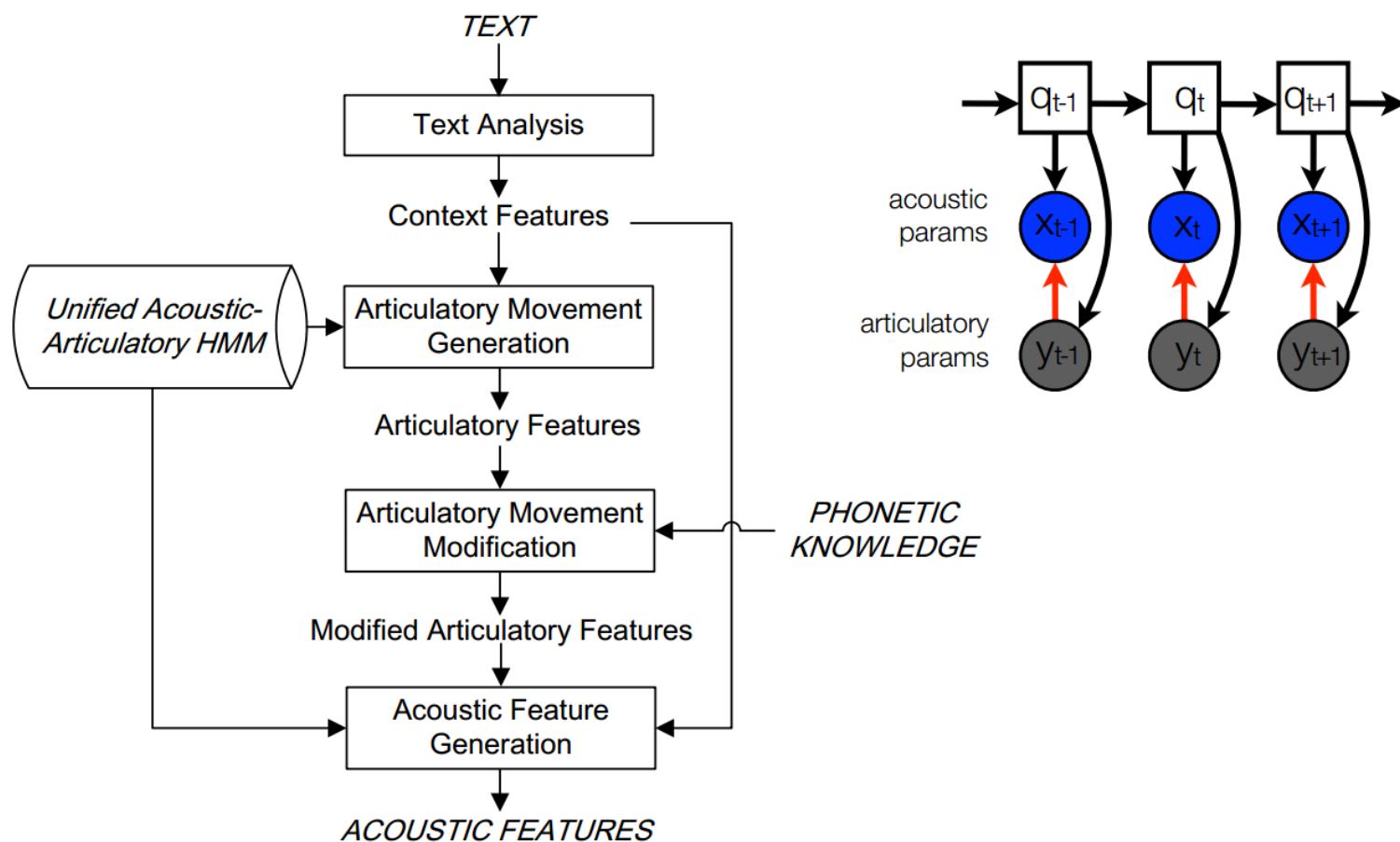  - Statistical processing
  - Computer controlled
- **Advantages**
  - Free vocabulary
  - Sound character modification

*SmartLab*
Intelligent Interactions
smartlab.tmit.bme.hu

BME TMIT

nVIDIA.
GPU
EDUCATION
CENTER

57/47

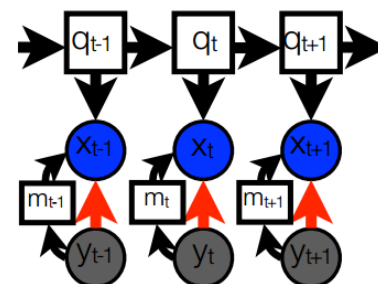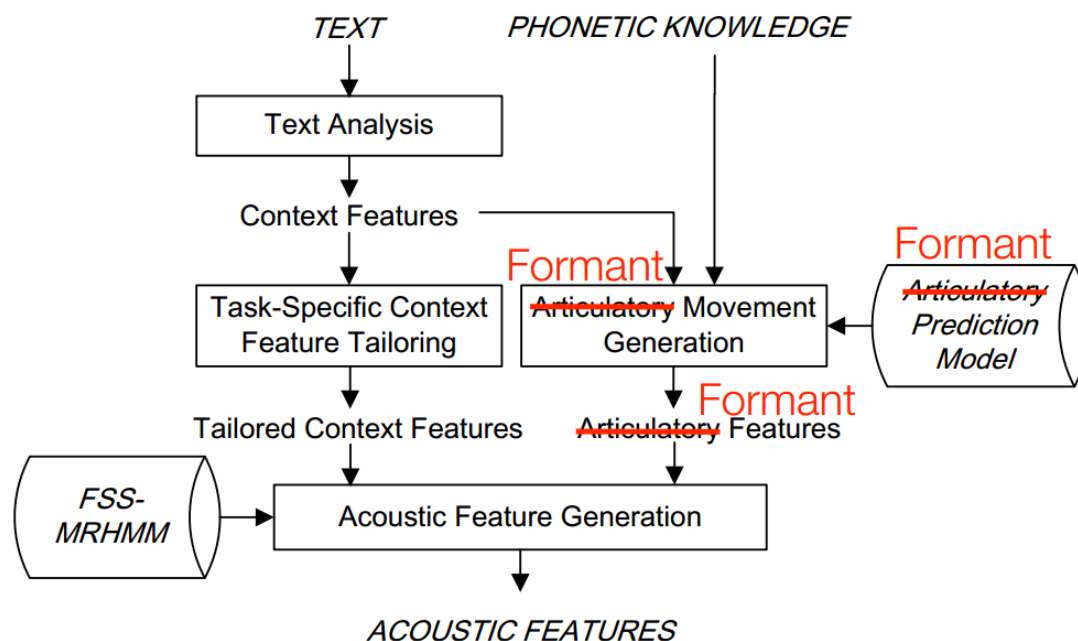# Articulatory/formant re-Ling et al . (2009-)



Synthesis with optional articulatory modification

# Articulatory/formant re-Ling et al . (2009-)

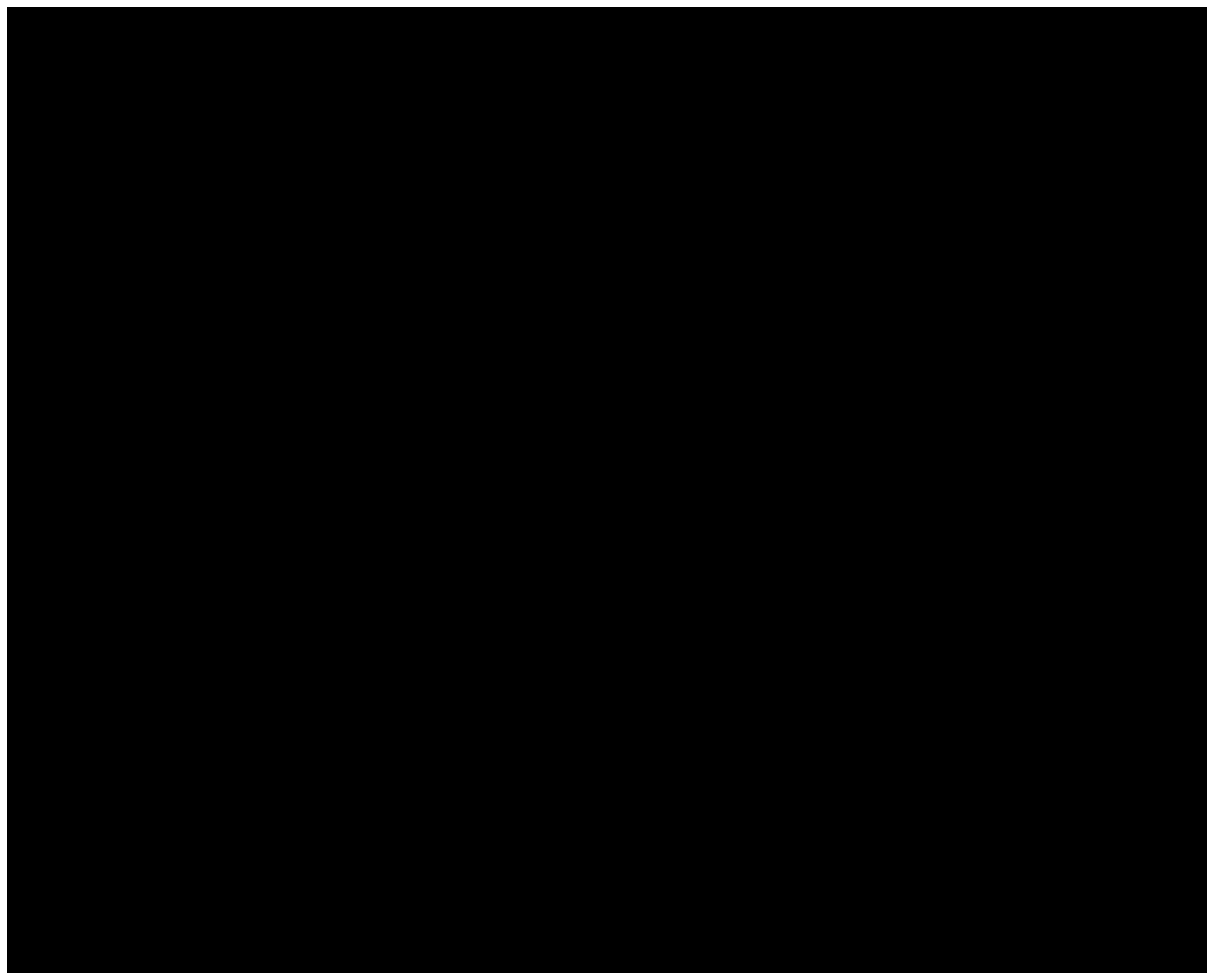## Alternatives to articulation - formants

# Articulatory synthesis
# Atsuo Takanishi (2008)



Digitally controlled mechanical articulation model

# Dyadic synthesis application (2012)
## Nao robot

# Outlook
## (even before deep learning)

- Challenges in machine speech generation
  - Adaptation
    - Application scenario (robot <> synthesized audiobook)
    - Speech rate (partner, message type)
    - Person and style (context, emotions, …)
    - Multilingualism (uniform quality, limited resources)
  - Intelligibility in noisy environments
    - Lombard effect
  - High-level control
    - Dialogue systems
  - . . . . . . . . . .

*SmartLab*
Intelligent Interactions
smartlab.tmit.bme.hu

BME TMIT

NVIDIA.
GPU
EDUCATION
CENTER

# Thank you for your attention!

**SmartLab**
Intelligent Interactions
smartlab.tmit.bme.hu

BME TMIT

NVIDIA.
GPU
EDUCATION
CENTER

63/47