

THE CURRENT SITUATION AND LONG-TERM PERSPECTIVES OF HUNGARIAN LANGUAGE TECHNOLOGY

Kornai András
MTA SZTAKI
BME Department of Algebra

CESAR

January 18 2013

WHEN DO WE CONSIDER HLT GOOD?

- When we don't notice it
- When it is error free
- When it becomes part of our everyday life
- When it expands our potential

WHEN DO WE CONSIDER HLT GOOD?

- When we don't notice it
 - When it is error free
 - When it becomes part of our everyday life
 - When it expands our potential

WHEN DO WE CONSIDER HLT GOOD?

- When we don't notice it
- When it is error free
- When it becomes part of our everyday life
- When it expands our potential

WHEN DO WE CONSIDER HLT GOOD?

- When we don't notice it
- When it is error free
- When it becomes part of our everyday life
- When it expands our potential

WHEN DO WE CONSIDER HLT GOOD?

- When we don't notice it
- When it is error free
- When it becomes part of our everyday life
- When it expands our potential

EXAMPLES

- Spellchecking
 - Not exactly unnoticeable
 - But indispensable
 - Not error free
 - But who dares to argue with the red underlining
- Web search
 - Not unnoticeable
 - Completely indispensable
 - Not error free
 - But the user doesn't know this

EXAMPLES

- Spellchecking

- 1 Not exactly unnoticeable
- 2 But indispensable
- 3 Not error free
- 4 But who dares to argue with the red underlining

- Web search

- 1 Not unnoticeable
- 2 Completely indispensable
- 3 Not error free
- 4 But the user doesn't know this

EXAMPLES

- Spellchecking
 - 1 Not exactly unnoticeable
 - 2 But indispensable
 - 3 Not error free
 - 4 But who dares to argue with the red underlining
- Web search
 - 1 Not unnoticeable
 - 2 Completely indispensable
 - 3 Not error free
 - 4 But the user doesn't know this

EXAMPLES

- Spellchecking
 - 1 Not exactly unnoticeable
 - 2 But indispensable
 - 3 Not error free
 - 4 But who dares to argue with the red underlining
- Web search
 - 1 Not unnoticeable
 - 2 Completely indispensable
 - 3 Not error free
 - 4 But the user doesn't know this

EXAMPLES

- Spellchecking
 - 1 Not exactly unnoticeable
 - 2 But indispensable
 - 3 Not error free
 - 4 But who dares to argue with the red underlining
- Web search
 - 1 Not unnoticeable
 - 2 Completely indispensable
 - 3 Not error free
 - 4 But the user doesn't know this

EXAMPLES

- Spellchecking
 - ① Not exactly unnoticeable
 - ② But indispensable
 - ③ Not error free
 - ④ But who dares to argue with the red underlining
- Web search
 - ① Not unnoticeable
 - ② Completely indispensable
 - ③ Not error free
 - ④ But the user doesn't know this

EXAMPLES

- Spellchecking
 - 1 Not exactly unnoticeable
 - 2 But indispensable
 - 3 Not error free
 - 4 But who dares to argue with the red underlining
- Web search
 - 1 Not unnoticeable
 - 2 Completely indispensable
 - 3 Not error free
 - 4 But the user doesn't know this

EXAMPLES

- Spellchecking
 - 1 Not exactly unnoticeable
 - 2 But indispensable
 - 3 Not error free
 - 4 But who dares to argue with the red underlining
- Web search
 - 1 Not unnoticeable
 - 2 Completely indispensable
 - 3 Not error free
 - 4 But the user doesn't know this

EXAMPLES

- Spellchecking
 - ① Not exactly unnoticeable
 - ② But indispensable
 - ③ Not error free
 - ④ But who dares to argue with the red underlining
- Web search
 - ① Not unnoticeable
 - ② Completely indispensable
 - ③ Not error free
 - ④ But the user doesn't know this

EXAMPLES

- Spellchecking
 - ① Not exactly unnoticeable
 - ② But indispensable
 - ③ Not error free
 - ④ But who dares to argue with the red underlining
- Web search
 - ① Not unnoticeable
 - ② Completely indispensable
 - ③ Not error free
 - ④ But the user doesn't know this

EXAMPLES

- Spellchecking
 - ① Not exactly unnoticeable
 - ② But indispensable
 - ③ Not error free
 - ④ But who dares to argue with the red underlining
- Web search
 - ① Not unnoticeable
 - ② Completely indispensable
 - ③ Not error free
 - ④ But the user doesn't know this

LANGUAGE TECHNOLOGIES IN DEVELOPMENT

Are very much like Machine Translation

- Not unnoticeable (but it would be if it were good)
- Not error free (but it would be if it were good)
- Not part of our everyday lives (but it would be if it were good)
- But it already expands our potential

LANGUAGE TECHNOLOGIES IN DEVELOPMENT

Are very much like Machine Translation

- Not unnoticeable (but it would be if it were good)
- Not error free (but it would be if it were good)
- Not part of our everyday lives (but it would be if it were good)
- But it already expands our potential

LANGUAGE TECHNOLOGIES IN DEVELOPMENT

Are very much like Machine Translation

- Not unnoticeable (but it would be if it were good)
- Not error free (but it would be if it were good)
- Not part of our everyday lives (but it would be if it were good)
- But it already expands our potential

LANGUAGE TECHNOLOGIES IN DEVELOPMENT

Are very much like Machine Translation

- Not unnoticeable (but it would be if it were good)
- Not error free (but it would be if it were good)
- Not part of our everyday lives (but it would be if it were good)
- But it already expands our potential

LANGUAGE TECHNOLOGIES IN DEVELOPMENT

Are very much like Machine Translation

- Not unnoticeable (but it would be if it were good)
- Not error free (but it would be if it were good)
- Not part of our everyday lives (but it would be if it were good)
- But it already expands our potential

WHAT DOES IT TAKE TO BUILD GOOD HLT?

- Largely correct theoretical model (e.g. Newtonian mechanics)
- Largely correct computational mechanism
- Largely correct/complete data
- Societal responsiveness

WHAT DOES IT TAKE TO BUILD GOOD HLT?

- Largely correct theoretical model (e.g. Newtonian mechanics)
- Largely correct computational mechanism
- Largely correct/complete data
- Societal responsiveness

WHAT DOES IT TAKE TO BUILD GOOD HLT?

- Largely correct theoretical model (e.g. Newtonian mechanics)
- Largely correct computational mechanism
- Largely correct/complete data
- Societal responsiveness

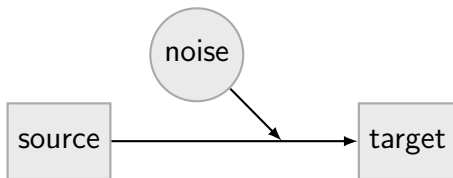
WHAT DOES IT TAKE TO BUILD GOOD HLT?

- Largely correct theoretical model (e.g. Newtonian mechanics)
- Largely correct computational mechanism
- Largely correct/complete data
- Societal responsiveness

WHAT DOES IT TAKE TO BUILD GOOD HLT?

- Largely correct theoretical model (e.g. Newtonian mechanics)
- Largely correct computational mechanism
- Largely correct/complete data
- Societal responsiveness

THE CLASSICAL FRAMEWORK (SHANNON-WEAVER 1948)



CHANGES SINCE 1948

- The classical framework was well suited for the traditional channels of communication **telephone, telex, printed matter**
- These channels were not capable of **processing, selecting, or transforming** information – these tasks were left entirely to the sender and the receiver.
- The chips in smartphones and in the computers that make up the internet have high computing power, which provides the foundation for **the channel itself** to provide a part of the intelligence.

CHANGES SINCE 1948

- The classical framework was well suited for the traditional channels of communication **telephone, telex, printed matter**
- These channels were not capable of **processing, selecting, or transforming** information – these tasks were left entirely to the sender and the receiver.
- The chips in smartphones and in the computers that make up the internet have high computing power, which provides the foundation for **the channel itself** to provide a part of the intelligence.

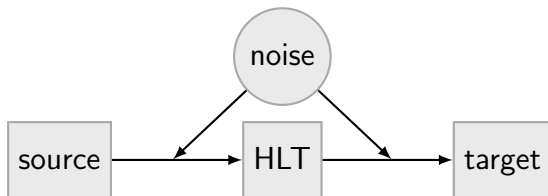
CHANGES SINCE 1948

- The classical framework was well suited for the traditional channels of communication **telephone, telex, printed matter**
- These channels were not capable of **processing, selecting, or transforming** information – these tasks were left entirely to the sender and the receiver.
- The chips in smartphones and in the computers that make up the internet have high computing power, which provides the foundation for **the channel itself** to provide a part of the intelligence.

CHANGES SINCE 1948

- The classical framework was well suited for the traditional channels of communication **telephone, telex, printed matter**
- These channels were not capable of **processing, selecting, or transforming** information – these tasks were left entirely to the sender and the receiver.
- The chips in smartphones and in the computers that make up the internet have high computing power, which provides the foundation for **the channel itself** to provide a part of the intelligence.

SHANNON-WEAVER 2.0



(HLT – Human Language Technology)

EXAMPLES

- Apple Spotlight – can find the form *went* if the query was *go*
- Browser can give cross-language pages (based e.g. on Google Translate)
- Microsoft speech translator

EXAMPLES

- Apple Spotlight – can find the form *went* if the query was *go*
- Browser can give cross-language pages (based e.g. on Google Translate)
- Microsoft speech translator

EXAMPLES

- Apple Spotlight – can find the form *went* if the query was *go*
- Browser can give cross-language pages (based e.g. on Google Translate)
- Microsoft speech translator

EXAMPLES

- Apple Spotlight – can find the form *went* if the query was *go*
- Browser can give cross-language pages (based e.g. on Google Translate)
- Microsoft speech translator

UNNOTICEABLE COMPONENTS

Like the carburetor: it is very hard to build a good one, it is indispensable for a working car, but the end-user is not at all interested

- Tokenization – where words and sentences begin and end
- Paer of speech tagging – nouns, verbs, adjectives, ...
- Morphological analysis – inflection, derivation, compounding, ...
- Named entity recognition – names of people, places, ...
- Chunking – finding Noun Phrases and other surface constituents
- Functional structure – subject, object, indirect object, ...
- Anaphors, pronoun resolution, dialog management, ...
- Logical inference, world model, ontology, ...
- ...

UNNOTICEABLE COMPONENTS

Like the carburetor: it is very hard to build a good one, it is indispensable for a working car, but the end-user is not at all interested

- Tokenization – where words and sentences begin and end
- Paer of speech tagging – nouns, verbs, adjectives, ...
- Morphological analysis – inflection, derivation, compounding, ...
- Named entity recognition – names of people, places, ...
- Chunking – finding Noun Phrases and other surface constituents
- Functional structure – subject, object, indirect object, ...
- Anaphors, pronoun resolution, dialog management, ...
- Logical inference, world model, ontology, ...
- ...

UNNOTICEABLE COMPONENTS

Like the carburetor: it is very hard to build a good one, it is indispensable for a working car, but the end-user is not at all interested

- Tokenization – where words and sentences begin and end
- Paer of speech tagging – nouns, verbs, adjectives, ...
- Morphological analysis – inflection, derivation, compounding, ...
- Named entity recognition – names of people, places, ...
- Chunking – finding Noun Phrases and other surface constituents
- Functional structure – subject, object, indirect object, ...
- Anaphors, pronoun resolution, dialog management, ...
- Logical inference, world model, ontology, ...
- ...

UNNOTICEABLE COMPONENTS

Like the carburetor: it is very hard to build a good one, it is indispensable for a working car, but the end-user is not at all interested

- Tokenization – where words and sentences begin and end
- Paer of speech tagging – nouns, verbs, adjectives, ...
- Morphological analysis – inflection, derivation, compounding, ...
- Named entity recognition – names of people, places, ...
- Chunking – finding Noun Phrases and other surface constituents
- Functional structure – subject, object, indirect object, ...
- Anaphors, pronoun resolution, dialog management, ...
- Logical inference, world model, ontology, ...
- ...

UNNOTICEABLE COMPONENTS

Like the carburetor: it is very hard to build a good one, it is indispensable for a working car, but the end-user is not at all interested

- Tokenization – where words and sentences begin and end
- Paer of speech tagging – nouns, verbs, adjectives, ...
- Morphological analysis – inflection, derivation, compounding, ...
- Named entity recognition – names of people, places, ...
- Chunking – finding Noun Phrases and other surface constituents
- Functional structure – subject, object, indirect object, ...
- Anaphors, pronoun resolution, dialog management, ...
- Logical inference, world model, ontology, ...
- ...

UNNOTICEABLE COMPONENTS

Like the carburetor: it is very hard to build a good one, it is indispensable for a working car, but the end-user is not at all interested

- Tokenization – where words and sentences begin and end
- Paer of speech tagging – nouns, verbs, adjectives, ...
- Morphological analysis – inflection, derivation, compounding, ...
- Named entity recognition – names of people, places, ...
- Chunking – finding Noun Phrases and other surface constituents
- Functional structure – subject, object, indirect object, ...
- Anaphors, pronoun resolution, dialog management, ...
- Logical inference, world model, ontology, ...
- ...

UNNOTICEABLE COMPONENTS

Like the carburetor: it is very hard to build a good one, it is indispensable for a working car, but the end-user is not at all interested

- Tokenization – where words and sentences begin and end
- Paer of speech tagging – nouns, verbs, adjectives, ...
- Morphological analysis – inflection, derivation, compounding, ...
- Named entity recognition – names of people, places, ...
- Chunking – finding Noun Phrases and other surface constituents
- Functional structure – subject, object, indirect object, ...
- Anaphors, pronoun resolution, dialog management, ...
- Logical inference, world model, ontology, ...
- ...

UNNOTICEABLE COMPONENTS

Like the carburetor: it is very hard to build a good one, it is indispensable for a working car, but the end-user is not at all interested

- Tokenization – where words and sentences begin and end
- Paer of speech tagging – nouns, verbs, adjectives, ...
- Morphological analysis – inflection, derivation, compounding, ...
- Named entity recognition – names of people, places, ...
- Chunking – finding Noun Phrases and other surface constituents
- Functional structure – subject, object, indirect object, ...
- Anaphors, pronoun resolution, dialog management, ...
- Logical inference, world model, ontology, ...
- ...

UNNOTICEABLE COMPONENTS

Like the carburetor: it is very hard to build a good one, it is indispensable for a working car, but the end-user is not at all interested

- Tokenization – where words and sentences begin and end
- Paer of speech tagging – nouns, verbs, adjectives, ...
- Morphological analysis – inflection, derivation, compounding, ...
- Named entity recognition – names of people, places, ...
- Chunking – finding Noun Phrases and other surface constituents
- Functional structure – subject, object, indirect object, ...
- Anaphors, pronoun resolution, dialog management, ...
- Logical inference, world model, ontology, ...
- ...

UNNOTICEABLE COMPONENTS

Like the carburetor: it is very hard to build a good one, it is indispensable for a working car, but the end-user is not at all interested

- Tokenization – where words and sentences begin and end
- Paer of speech tagging – nouns, verbs, adjectives, ...
- Morphological analysis – inflection, derivation, compounding, ...
- Named entity recognition – names of people, places, ...
- Chunking – finding Noun Phrases and other surface constituents
- Functional structure – subject, object, indirect object, ...
- Anaphors, pronoun resolution, dialog management, ...
- Logical inference, world model, ontology, ...
- ...

HUNGARIAN LANGUAGE TECHNOLOGY BASED ON ENGLISH MODELS?

- Feasible for speech, but the English model **can't take advantage of** the major regularities of Hungarian speech such as vowel harmony
- Not feasible at the word level, since the English model **can't cope with** the complex Hungarian patterns of suffixation
- Feasible for analyzing NPs, since the English model **must be prepared for about the same complexity** as the Hungarian model (fixed word order within NPs)
- Not feasible at the sentence level, since the English model **is based on a condition not met by Hungarian** (fixed word order)

HUNGARIAN LANGUAGE TECHNOLOGY BASED ON ENGLISH MODELS?

- Feasible for speech, but the English model **can't take advantage of** the major regularities of Hungarian speech such as vowel harmony
- Not feasible at the word level, since the English model **can't cope with** the complex Hungarian patterns of suffixation
- Feasible for analyzing NPs, since the English model **must be prepared for about the same complexity** as the Hungarian model (fixed word order within NPs)
- Not feasible at the sentence level, since the English model **is based on a condition not met by Hungarian** (fixed word order)

HUNGARIAN LANGUAGE TECHNOLOGY BASED ON ENGLISH MODELS?

- Feasible for speech, but the English model **can't take advantage of** the major regularities of Hungarian speech such as vowel harmony
- Not feasible at the word level, since the English model **can't cope with** the complex Hungarian patterns of suffixation
- Feasible for analyzing NPs, since the English model **must be prepared for about the same complexity** as the Hungarian model (fixed word order within NPs)
- Not feasible at the sentence level, since the English model **is based on a condition not met by Hungarian** (fixed word order)

HUNGARIAN LANGUAGE TECHNOLOGY BASED ON ENGLISH MODELS?

- Feasible for speech, but the English model **can't take advantage of** the major regularities of Hungarian speech such as vowel harmony
- Not feasible at the word level, since the English model **can't cope with** the complex Hungarian patterns of suffixation
- Feasible for analyzing NPs, since the English model **must be prepared for about the same complexity** as the Hungarian model (fixed word order within NPs)
- Not feasible at the sentence level, since the English model **is based on a condition not met by Hungarian** (fixed word order)

WHAT DOES HUNGARIAN HLT REQUIRE?

- Largely correct theoretical model
- Largely correct computational mechanism
- Largely correct/complete data
- Societal responsiveness

WHAT DOES HUNGARIAN HLT REQUIRE?

- Largely correct theoretical model
- Largely correct computational mechanism
- Largely correct/complete data
- Societal responsiveness

WHAT DOES HUNGARIAN HLT REQUIRE?

- Largely correct theoretical model
- Largely correct computational mechanism
- Largely correct/complete data
- Societal responsiveness

WHAT DOES HUNGARIAN HLT REQUIRE?

- Largely correct theoretical model
- Largely correct computational mechanism
- Largely correct/complete data
- Societal responsiveness

WHAT DOES HUNGARIAN HLT REQUIRE?

- Largely correct theoretical model
- Largely correct computational mechanism
- Largely correct/complete data
- Societal responsiveness

WHAT IS THE STATE OF TEH ART?

- Largely correct theoretical model – still requires a lot of work
- Largely correct computational mechanism – we are doing well
- Largely correct/complete data – we are not doing too well
- Societal responsiveness – not without problems

WHAT IS THE STATE OF THE ART?

- Largely correct theoretical model – still requires a lot of work
- Largely correct computational mechanism – we are doing well
- Largely correct/complete data – we are not doing too well
- Societal responsiveness – not without problems

WHAT IS THE STATE OF THE ART?

- Largely correct theoretical model – still requires a lot of work
- Largely correct computational mechanism – we are doing well
- Largely correct/complete data – we are not doing too well
- Societal responsiveness – not without problems

WHAT IS THE STATE OF THE ART?

- Largely correct theoretical model – still requires a lot of work
- Largely correct computational mechanism – we are doing well
- Largely correct/complete data – we are not doing too well
- Societal responsiveness – not without problems

WHAT IS THE STATE OF THE ART?

- Largely correct theoretical model – still requires a lot of work
- Largely correct computational mechanism – we are doing well
- Largely correct/complete data – we are not doing too well
- Societal responsiveness – not without problems

MAJOR RESEARCH CENTERS

- Hungarian Academy of Sciences (Institute of Linguistics, Computer and Automation Research Institute, ...)
- Universities (BUTE, Debrecen, PPCU, Szeged,...)
- Hungarian industry (Aitia, ALL, Kilgray, Morphologic,...)
- Hungarian products by multinationals (Apple, Google, IBM, Microsoft, Nuance, Xerox,...)

MAJOR RESEARCH CENTERS

- Hungarian Academy of Sciences (Institute of Linguistics, Computer and Automation Research Institute, ...)
- Universities (BUTE, Debrecen, PPCU, Szeged,...)
- Hungarian industry (Aitia, ALL, Kilgray, Morphologic,...)
- Hungarian products by multinationals (Apple, Google, IBM, Microsoft, Nuance, Xerox,...)

MAJOR RESEARCH CENTERS

- Hungarian Academy of Sciences (Institute of Linguistics, Computer and Automation Research Institute, ...)
- Universities (BUTE, Debrecen, PPCU, Szeged,...)
- Hungarian industry (Aitia, ALL, Kilgray, Morphologic,...)
- Hungarian products by multinationals (Apple, Google, IBM, Microsoft, Nuance, Xerox,...)

MAJOR RESEARCH CENTERS

- Hungarian Academy of Sciences (Institute of Linguistics, Computer and Automation Research Institute, ...)
- Universities (BUTE, Debrecen, PPCU, Szeged,...)
- Hungarian industry (Aitia, ALL, Kilgray, Morphologic,...)
- Hungarian products by multinationals (Apple, Google, IBM, Microsoft, Nuance, Xerox,...)

MAJOR RESEARCH CENTERS

- Hungarian Academy of Sciences (Institute of Linguistics, Computer and Automation Research Institute, ...)
- Universities (BUTE, Debrecen, PPCU, Szeged,...)
- Hungarian industry (Aitia, ALL, Kilgray, Morphologic,...)
- Hungarian products by multinationals (Apple, Google, IBM, Microsoft, Nuance, Xerox,...)

Thank you for your attention!

SOCIETAL RESPONSIVENESS

An article from a leading Hungarian technology blog about SZTAKI's language-driven railroad ticket clerk
[Click through the above line only if you read Hungarian]

SOCIETAL RESPONSIVENESS II

Erika | 2012. 09. 28. 21:06

2

Tehát az eddig is közutálatnak örvendő hangosmenü géphangot átveszi a szoftver... ez biztos sokkal közkedveltebb és népszerűbb lesz, kérdés csupán az, hogy mitől???

caramell | 2012. 09. 28. 20:17

1

"Kornai András szerint a megoldásuk nem csupán a MÁV-pénztáros szerepét tudná átvenni. A kutató szerint az összes olyan emberi munkát ki lehetne váltani szoftverrel, ahová csak azért ültetnek embert, mert az képes mondatokat értelmezni. Még akár a telefonos ügyfélszolgálatok közutálatnak örvendő hangosmenüjének helyét is átveheti egy olyan szoftver, aminek el kell mondani, hogy milyen panasszal telefonáltak, és rögtön kapcsolja a megfelelő ügyintézős."

HATÁROZOTTAN ELLENE VAGYOK,IDEJE LENNE HATÁRT SZABNI AZ EMBEREK FÖLDÖNFUTÓVÁ TÉTELÉNEKINEM ERRE KELLENE HASZNÁLNI AZ INFORMATIKÁT MOST ILYEN EGYÉBKÉNT IS NEHÉZ GAZDASÁGI VÁLSÁGBAN!

EZZEL NEM MUNKAHELYET VÉDENEKI,HANEM MEGSEMISITENEK.

CSAK A PROFIT SZÁMIT!? AZ Ember,A CSALÁDOK FÖLDÖNFUTÓVÁ TÉTELE NEM?
EZZEL CSAK AZ 1-2 EGYÉBKÉNT IS MILLIÁRDOS INFORMATIKUST TESZIK MÉG GAZDAGABBÁ,MIKÖZBEN AZ EMBEREK EZREIT,VAGY MILLIÓIT MUNKANÉLKÜLIVÉ TESZNEK.

[Reader comment on the article, accusing 'billionaire informatics experts' of destroying jobs]

SOCIETAL RESPONSIVENESS III

[origo]

EZT OLVASTA MÁR?

◀ 1 / 5 ▶



Ijesztően sokat tud a Facebook új keresője

[Recommendation from the same tech blog – Facebook's new search 'knows a frightening amount']