**FACULTY
OF INFORMATICS**
Masaryk University

# Trains and Trees of Thoughts

## Towards the Representation of Structural Semantics

**Petr Sojka**

`sojka@fi.muni.cz`

Semiomaths workshop, ETH Zurich, 2018, March 2nd, 2PM

## Brno, Czech Republic

■ town, where Kurt Gödel was born
and spent his childhood

■ town, where Gregor Mendel has founded genetics

■ town of 'Silicon Valley' of [Central] Europe with high concentration of
Computer Science bussinesses (RedHat, IBM, Kiwi, Honeywell), and
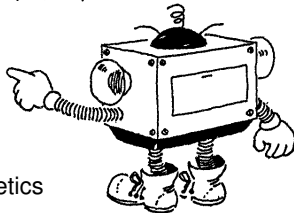academia (Faculty of Informatics MU, FIT MU)

# Table of Contents

# Semiotics

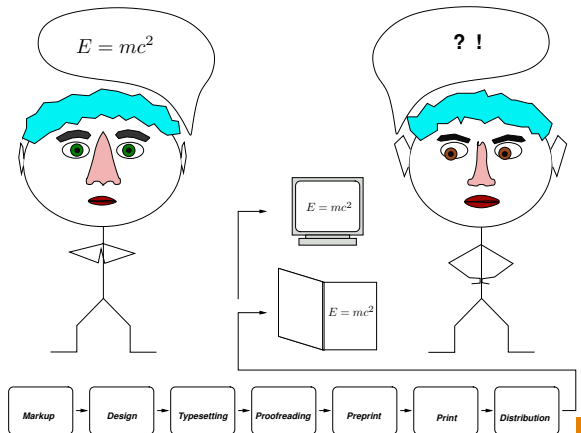Study of signs as means of language or *communication*.

Is math communication specific?

Are the signs used in math specific?

Are computer programming languages specific?

How mathematics and computer science differ?

# Scholarly Communication via Digital Mathematics Libraries (DMLs): DML-CZ, EuDML project participation

**FACULTY**
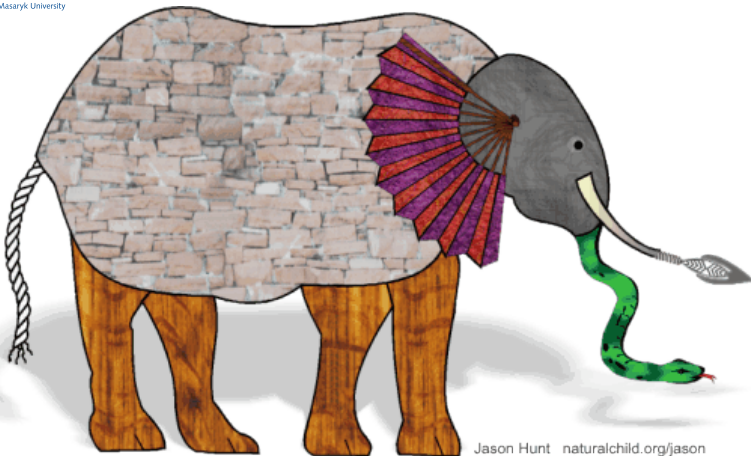**OF INFORMATICS**
Masaryk University

# Paul Watzlawick's Five Axioms of Communication

1. (cannot not) One cannot not communicate.

2. (content & relationship) Every communication has a content and relationship aspect such that the latter classifies the former and is therefore a meta-communication.

3. (punctuation) The nature of a relationship is dependent on the punctuation of the partners' communication procedures.

4. (digital & analogic, discrete & continuous) Human communication involves both digital and analogic modalities.

5. (symmetric or complementary) Inter-human communication procedures are either symmetric or complementary, depending on whether the relationship of the partners is based on differences or parity.
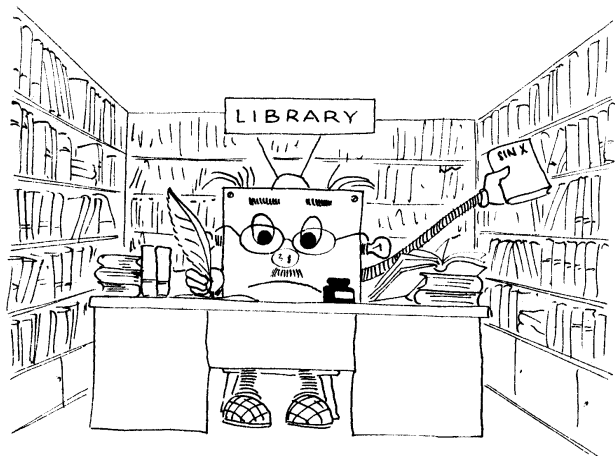
Do they hold for math-specific discourse?

# Six blind monks examining an elephant on a trip

Jason Hunt   naturalchild.org/jason

Q: Is elephant a wall (belly), hand fan (ear), solid pipe (tusk), pillar (leg), rope (tail) or tree branch (trunk)?

Let the animal on the road is meaning-conveying [math] communication in the form of scientific papers stored in the digital libraries like the EuDML or arXiv, digital library with math content.

# Six blind "monks" quarreling after the trip

1. (Shannon) touched information-theoretic properties of communication

2. (Chomsky) stresses formal grammars in communication languages

3. (Gödel) states limitations in formal expressiveness communication, being incomplete or inconsistent

4. (Rogers) person-centered communication

5. (Locke) views that knowledge comes primarily from sensory experience: empiricism

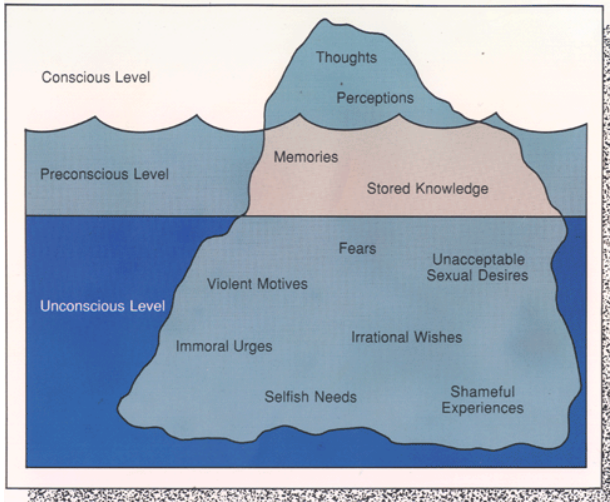6. (Lakoff) argues that conceptual metaphors are basis for embodied minds' communication

Who is right? Does the question makes sense at all? Whom do you identify with on your ride of thoughts?

**FACULTY
OF INFORMATICS**
Masaryk University

# Metaphor of an Iceberg

- Icebergs have small visible and big invisible part.

- Icebergs with the same visible parts may have quite different invisible parts.

- Iceberg does have internal structure and qualities that differ.

- The glaciers swim, struggle against each other, melt down, diminish, grow, …

***

- Signs are created in human minds: they reflect human's (objective) conscious thoughts, but also human's (subjective) unconscious personality.

- Freud's view of the human mind is *mental iceberg*.

**FACULTY
OF INFORMATICS**
Masaryk University

# The Mental Iceberg metaphor: conscious, preconscious and unconsious levels

Only 10% of an iceberg is visible (conscious) whereas the other 90% is beneath the water: the preconscious is allotted approximately 10%–15% whereas the unconscious is allotted an overwhelming 75%–80%.

Conscious $\equiv$ visible (surface) words/texts, same for all.

Preconscious $\equiv$ structured information deducible from (surface) text based on language and common knowledge, personal (e.g. different and subjective based on previous occurences). Mostly present *latently*.

Unconscious $\equiv$ yet unnamed unknown relations and knowledge indirectly related to the meaning of the conveyed, visible part of mind: expressed messages.

FACULTY
OF INFORMATICS
Masaryk University

# The Semiotics Iceberg Metaphor: layers of *visible* signs and hidden *shared* and *personal* structures linked to them

Only 10% of communication is on a) **surface** (written text with visual marks and punctuation) whereas the other 90% is beneath the water: 'preconscious' b) **shared common** sense layer, and 'subconscious' level of c) **personal** association and connotations. agreed notions (arithmetics) and notation (formulas) ground preconscious and unconscious).

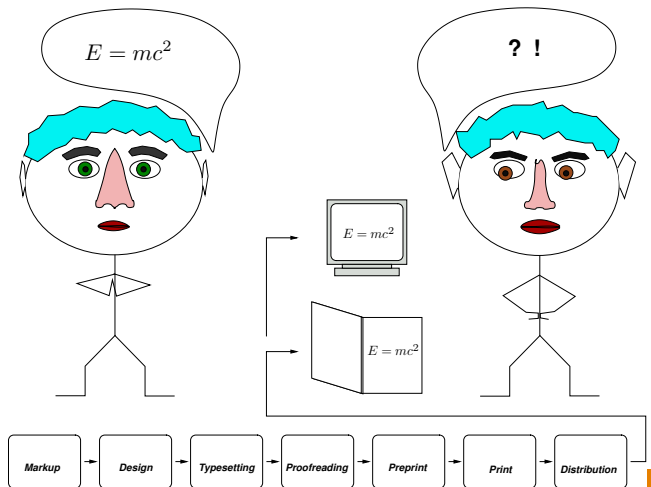a) $\equiv$ visible (surface) words/texts, formulas.

b) $\equiv$ structured information deducible from (surface) text based on language and common knowledge, latently present in the communication. Mostly present *latently*.

c) $\equiv$ subjective, possibly emotional connotation the conveyed message raises.

**FACULTY
OF INFORMATICS**
Masaryk University

# The Semiotics Iceberg Metaphor: $E = mc^2$

- *visible* part: $E = mc^2$

- signs and hidden *shared* part: notions of mass, energy, speed of light and their notations

- *personal* structures and connotations linked to them

# Scholarly Communication *via DMLs using rich KB*

# Sketch Engine: from statistics to insight

**FACULTY
OF INFORMATICS**
Masaryk University

# Which insight math corpora could give us? Pros and cons

Firthian linguistics: You shall know a word by the company it keeps (Firth, J.R., 1957).

You shall know a sign/ notion/ math formulae by the company it keeps.

- Sketches of math signs similarly generated as word sketches?

- Math corpora reveals globality or locality of sign/ notion usage.

- Putting knowledge on one place allows for new killer application: search, similar sign search, similar phrase search, similar formulae search, similar plagiarism search based on word n-grams similar theorem search similar thoughts and structures search similar XXX search.

# Towards higher level content representations – knowledge bases

NLP processing from strings via words to meaning, including *math-awareness*
math specifics: structures and abstractions

- to allow searching (semantically) similar papers, precise [semantic] indexing: search as a gate to knowledge

- to allow exploration of a DML by intelligent browsing of (semantically) man is known by the company he keepssimilar papers: distributional semantics topic modeling as Latent Semantic Indexing, Latent Dirichlet Allocation

- to allow personalization and domain specifics, e.g. semantic faceted search (formulae,…)

- to track 'train of thought' – narrative qualities of papers, proofs (Mizar type of paper)

FACULTY
OF INFORMATICS
Masaryk University

# Motivation for example I

```
From: Shayan A Tabrizi <shayantabrizi@gmail.com>
Subject: [Corpora-List] Dataset for Different Research Areas

I want to find the relevance of each of the research papers of
my dataset to each of the research areas such as Physics, CS,
Math, Social Sciences, etc.
Thus, I need a dataset consisting of all research areas and
some sample texts (preferably papers) in that area, to estimate
the similarity of each of my papers to each of the areas.
               *Is there any such dataset?*
Some points:
   1. It is much much better if the dataset has areas in different
   granularities. e.g. in one level: Mathematics, Physics, CS, etc.
   and in a more fine-grained level divides CS to Networks,
   Artificial Intelligence, etc.
   2. Even if the dataset only consists of a specific domain
   (especially CS) and its sub-domains it is still usable.
```

# Probabilistic Topical Modeling: Latent Dirichlet Allocation

- topic: weighted list of words

- document: weighted list of topics

# Topical Modeling: Latent Dirichlet Allocation II

■ all topics computed automatically from document corpora



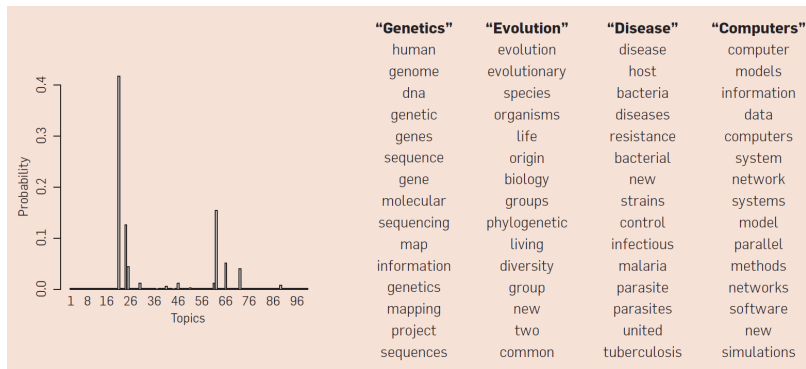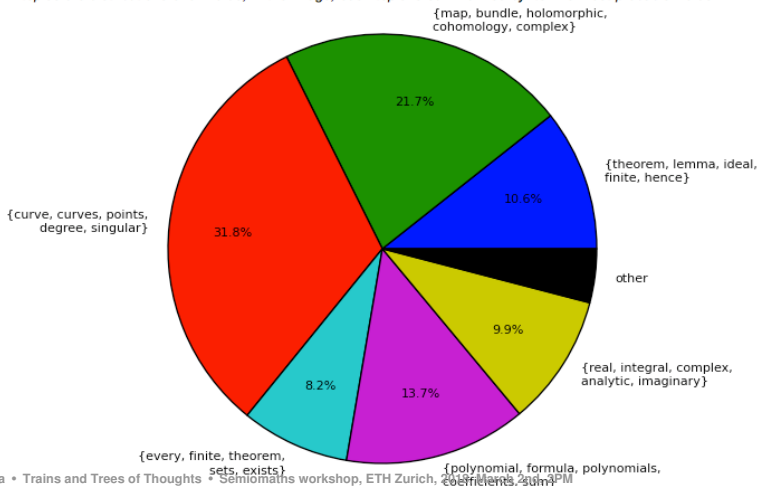| | "Genetics" | "Evolution" | "Disease" | "Computers" |
|---|---|---|---|---|
| | human | evolution | disease | computer |
| | genome | evolutionary | host | models |
| | dna | species | bacteria | information |
| | genetic | organisms | diseases | data |
| | genes | life | resistance | computers |
| | sequence | origin | bacterial | system |
| | gene | biology | new | network |
| | molecular | groups | strains | systems |
| | sequencing | phylogenetic | control | model |
| | map | living | infectious | parallel |
| | information | diversity | malaria | methods |
| | genetics | group | parasite | networks |
| | mapping | new | parasites | software |
| | project | two | united | new |
| | sequences | common | tuberculosis | simulations |

# Example I: Automated Meaning Picking from Texts



**LDA** Topics Pie Chart for **math.0406240**:
*Each slice represents a different topic. The size of the slice corresponds to "how much is the article about this topic?". Topics which contribute <6% to the above document are aggregated under "other".*
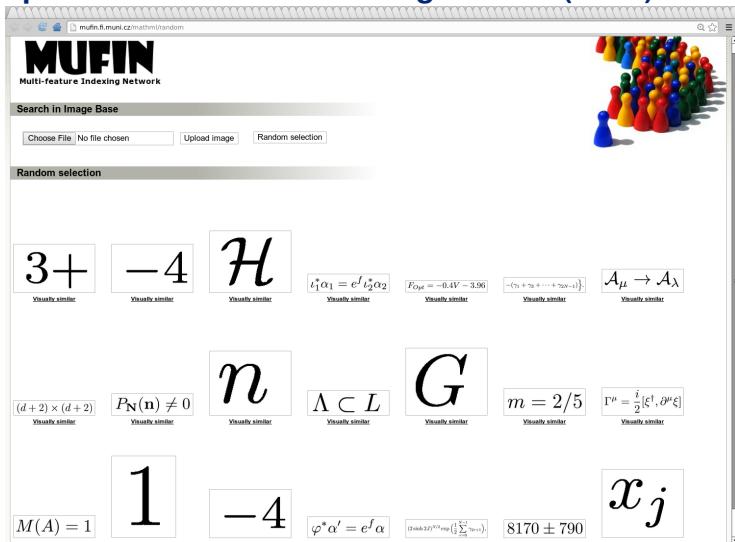
*LDA topics are distributions over words; in the image, each topic is summarized by its five most probable words.*

# Picking characteristic words from 'top of iceberg': LDA topics for given arXiv article

<https://mir.fi.muni.cz/eudmldemo/gensim-arxiv/>

# Example II: metric-based indexing of math (+text)

# Example III: text parsing with ParsCit

From OCR we get:

```
[5] Lambe, L., Stasheff, J.: Applications
of perturbation theory to iterated fibrations.
Manuscripta Math. 58 (1987), 363–376.
```

**FACULTY
OF INFORMATICS**
Masaryk University

## Parsing citations with ParsCit

```
<algorithms version="110505">
  <algorithm name="ParsCit" version="110505">
    <citationList>
      <citation valid="true">
        <authors>
          <author>L Lambe</author>
          <author>J Stasheff</author>
        </authors>
        <title>Applications of perturbation theory to iterated
                fibrations.</title>
        <date>1987</date>
        <journal>Manuscripta Math.</journal>
        <volume>58</volume>
        <pages>363--376</pages>
        <marker>[5]</marker>
        <rawString>Lambe, L., Stasheff, J.: Applications of
                   perturbation theory to iterated fibrations.
                   Manuscripta Math. 58 (1987), 363-376.</rawString>
      </citation>
    </citationList>
  </algorithm>
</algorithms>
```

# Word representations: Word2vec

- Tomas Mikolov from Brno came with the idea of machine learnt representation of words in high-dimensional spaces.

- The representation *capture* both syntactic and semantic properties of word usage in context.

- Only global properties are represented (not outliers).

- This continuous representation proved superior to previous discrete representation of words (of Wordnet type).

FACULTY
OF INFORMATICS
Masaryk University

# Representational learning

"Le silence eternel de ces espaces infinis m'effraie."

"Those eternal silence of these infinite spaces terrifies me."

Blaise Pascal, 1670

- representations in different spaces (vector ones, hyperbolic ones …)

- curse of dimensionality, projections, serialization

- ambiguity raises in low dimensions (is it practical or not?)

# Learning hierarchical representations: adding structural qualitites of communicated texts

- Poincare embeddings (Nickel, Kiela, 2017) could learn latent hierarchical structural qualities (collocations, phrases, formulae trees, …).

- Trajectories of these representations (points on an $n$-dimensional Poincaré ball) could represent complex hierarchical entities 'structural signs': thoughts.

# Table of Contents

FACULTY
OF INFORMATICS
Masaryk University

# Six blind monks examining an elephant on a trip

FACULTY
OF INFORMATICS
Masaryk University

# Questions?

# Acknowledgements

Organizers of Semiomaths workshop for invitation.

# References

[1]    Archambault, D., Moço, V.: Canonical MathML to Simplify Conversion of MathML to Braille
       Mathematical Notations. In: Miesenberger, K., Klaus, J., Zagler, W., Karshmer, A. (eds.)
       Computers Helping People with Special Needs, Lecture Notes in Computer Science, vol.
       4061, pp. 1191–1198. Springer Berlin / Heidelberg (2006),
       <https://doi.org/10.1007/11788713_172>

[2]    Grimm, J.: Producing MathML with Tralics. In: Sojka [4], pp. 105–117,
       <https://dml.cz/dmlcz/702579>

[3]    MREC – Mathematical REtrieval Collection, <https://mir.fi.muni.cz/MREC/index.html>

[4]    Sojka, P. (ed.): Towards a Digital Mathematics Library. Masaryk University, Paris, France
       (Jul 2010), <https://www.fi.muni.cz/ sojka/dml-2010-program.html>

[5]    Sojka, P. (ed.): Towards a Digital Mathematics Library. Masaryk University, Paris, France
       (Jul 2011), <https://www.fi.muni.cz/ sojka/dml-2011-program.html>

# References (cont.)

[6]    Sojka, P., Líška, M.: Indexing and Searching Mathematics in Digital Libraries – Architecture,
       Design and Scalability Issues. In: Davenport, J.H., Farmer, W., Rabe, F., Urban, J. (eds.)
       Proceedings of CICM Conference 2011 (Calculemus/MKM). Lecture Notes in Artificial
       Intelligence, LNAI, vol. 6824, pp. 228–243. Springer-Verlag, Berlin, Germany (Jul 2011)

[7]    Stamerjohanns, H., Ginev, D., David, C., Misev, D., Zamdzhiev, V., Kohlhase, M.:
       MathML-aware Article Conversion from LaTeX. In: Sojka, P. (ed.) Proceedings of DML 2009.
       pp. 109–120. Masaryk University, Grand Bend, Ontario, CA (Jul 2009),
       <https://dml.cz/dmlcz/702561>

[8]    Stamerjohanns, H., Kohlhase, M., Ginev, D., David, C., Miller, B.: Transforming Large
       Collections of Scientific Publications to XML. Mathematics in Computer Science 3, 299–307
       (2010), <https://doi.org/10.1007/s11786-010-0024-7>

[9]    Sylwestrzak, W., Borbinha, J., Bouche, T., Nowiński, A., Sojka, P.: EuDML—Towards the
       European Digital Mathematics Library. In: Sojka [4], pp. 11–24,
       <https://dml.cz/dmlcz/702569>

# References (cont.)

[10] Líška, Martin and Petr Sojka and Michal Růžička. Similarity Search for Mathematics: Masaryk University team at the NTCIR-10 Math T ask. In Proceedings of the 10th NTCIR Workshop Meeting on Evaluation of Information Access Technologies: Math Pilot Task. pp. 686-691. NII, Tokyo, 2013. PDF

[11] D. Formánek, M. Líška, M. Růžička, and P. Sojka. Normalization of digital mathematics library content. In J. Davenport, J. Jeuring, C. Lange, and P. Libbrecht, editors, 24th OpenMath Workshop, 7th Workshop on Mathematical User Interfaces (MathUI), and Intelligent Computer Mathematics Work in Progress, number 921 in CEUR Workshop Proceedings, pp. 91–103, Aachen, 2012.

[12] Sojka, Petr and Martin Líška. The Art of Mathematics Retrieval. In Matthew R. B. Hardy , Frank Wm. Tompa. Proceedings of the 2011 ACM Symposium on Document Engineering. Mountain View, CA, USA: ACM, 2011. p. 57–60. ISBN 978-1-4503-0863-2. <http://dx.doi.org/10.1145/2034691.2034703>

# References (cont.)

[13]   Martin Líška, Petr Sojka, Michal Růžička, and Petr Mravec.
       Web Interface and Collection for Mathematical Retrieval.
       In Petr Sojka and Thierry Bouche, editors, *Proceedings of DML 2011*, pages 77–84,
       Bertinoro, Italy, July 2011. Masaryk University.
       <https://dml.cz/dmlcz/702604>.

[14]   Credits for LDA pictures goes to David M. Blei.

[15]   Credits for illustrations goes to Jiří Franek.