

# Understanding Medical Data: Text Analysis and Coding with Semfinder Expert System

Jörg Niggemann, Michael Owsijewitsch, Hans Rudolf Straub  
Semfinder AG, Kreuzlingen, Switzerland [www.semfinder.com](http://www.semfinder.com) [straub@semfinder.com](mailto:straub@semfinder.com)



**Challenge** Unstructured information represents the largest and most relevant source of information in hospitals

**The Result** Software solution in daily use for coding:

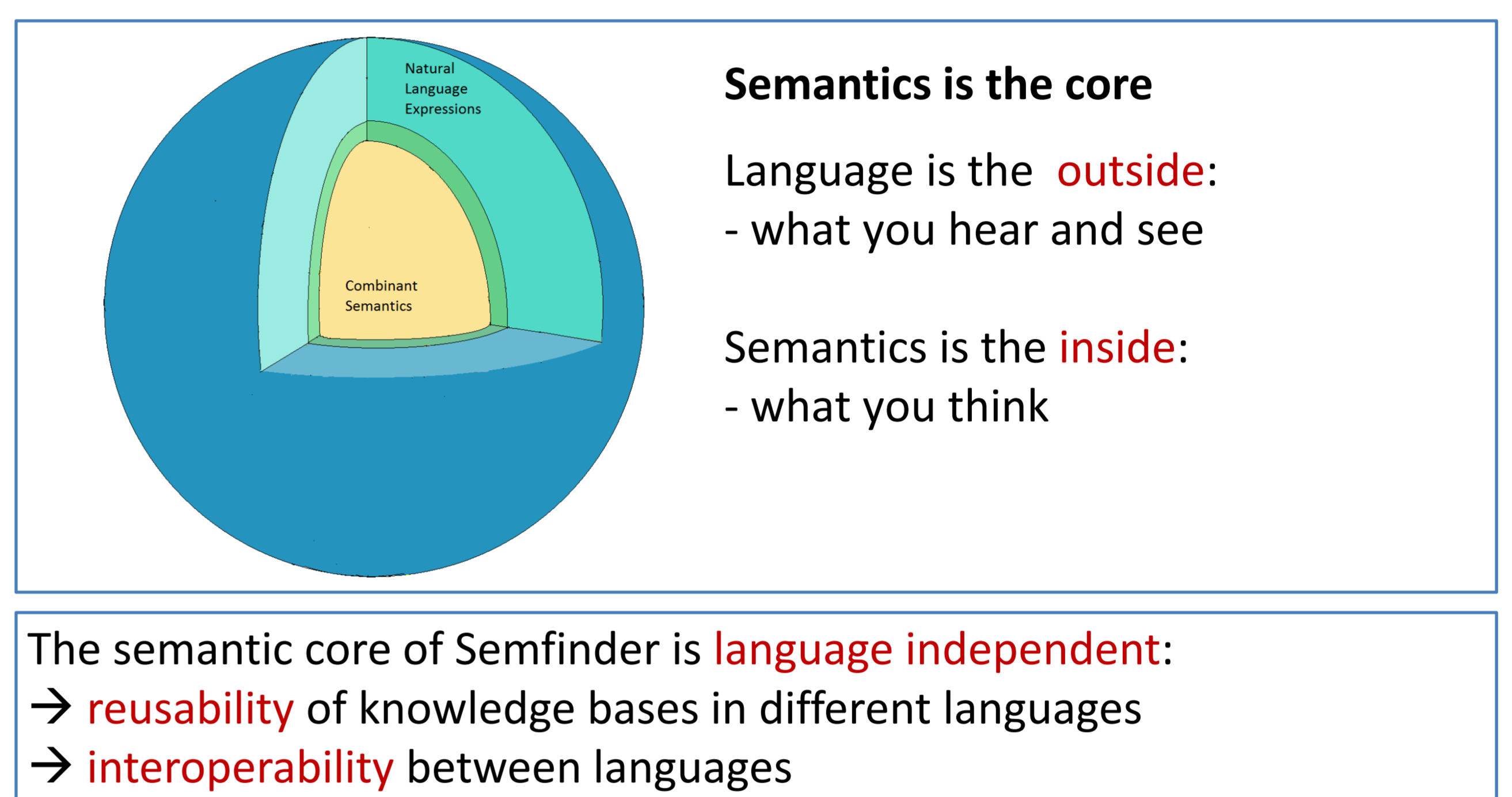
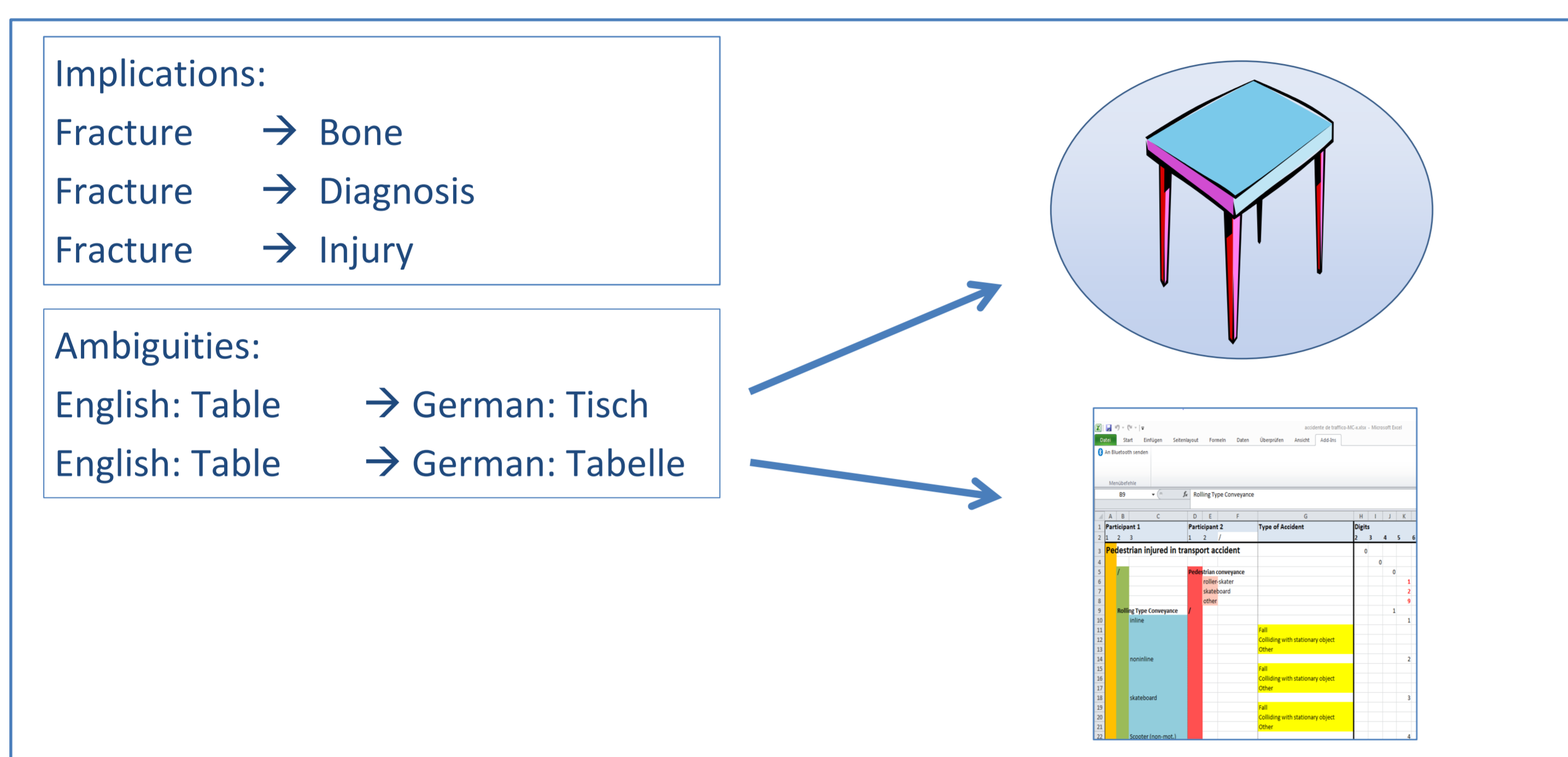
- 500+ hospitals in Germany, Switzerland and Spain
- Integration in IBM Insurance Service Hub for analysis of 2 Mio diagnoses / month in Germany

## Sources of complexity in medical free texts

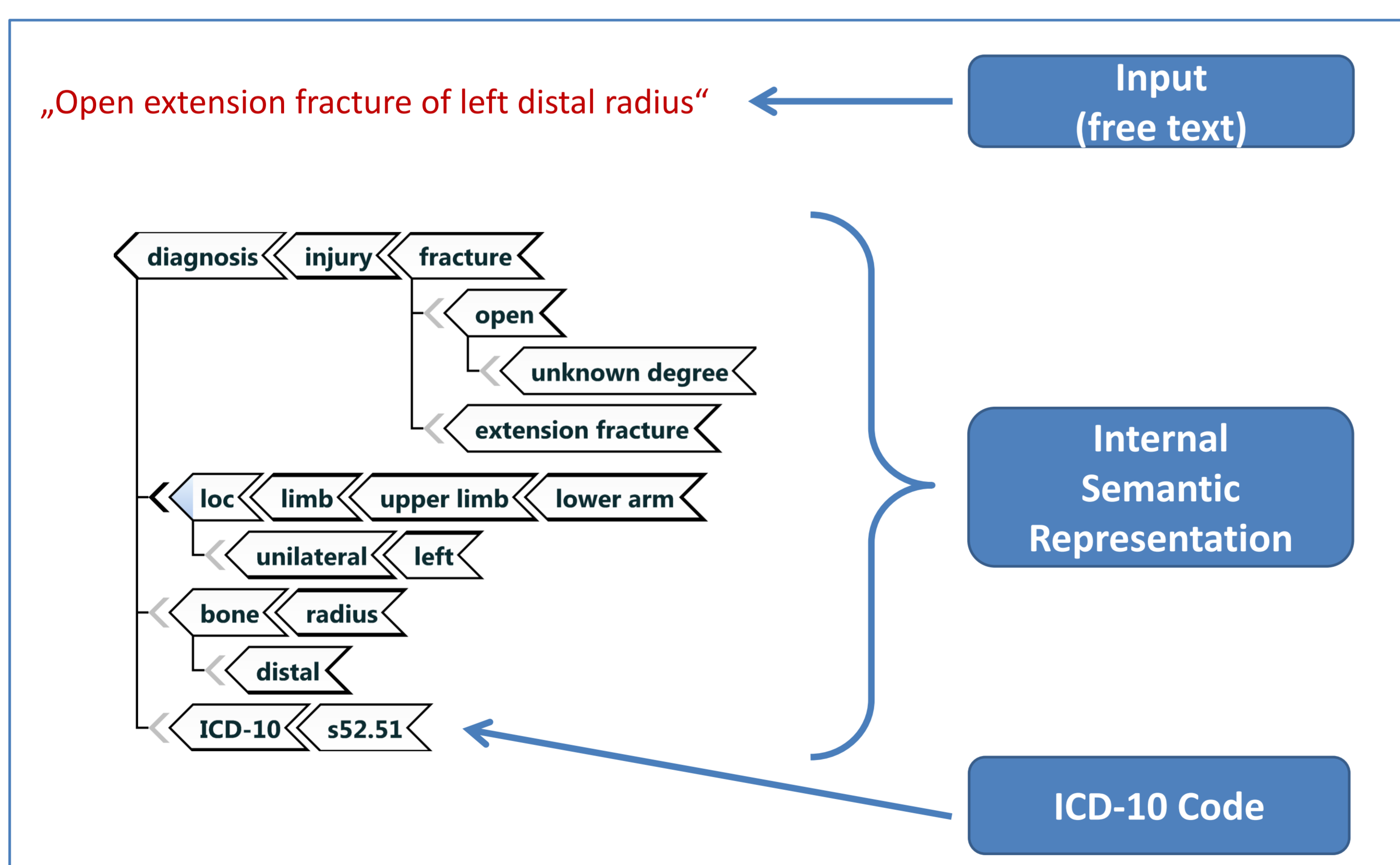
- Overlappings (*streptococcal pneumonia, postoperative*)
- Ambiguities (*heads in the shoulders and abdomen*)
- Negations (*diabetes, non-insulin-dependent, with complications*)
- Non-information (*diabetes, not otherwise specified*)
- Implications (*radius → bone, forearm*)
- Omissions (*fracture of humerus and radius*)
- Composite Diagnoses with mutual dependencies

## Semantics: words ≠ concepts

Typical for semantics: → **Overlapping** of concepts  
→ **Implicit** concepts  
→ Resolvable **Ambiguities**



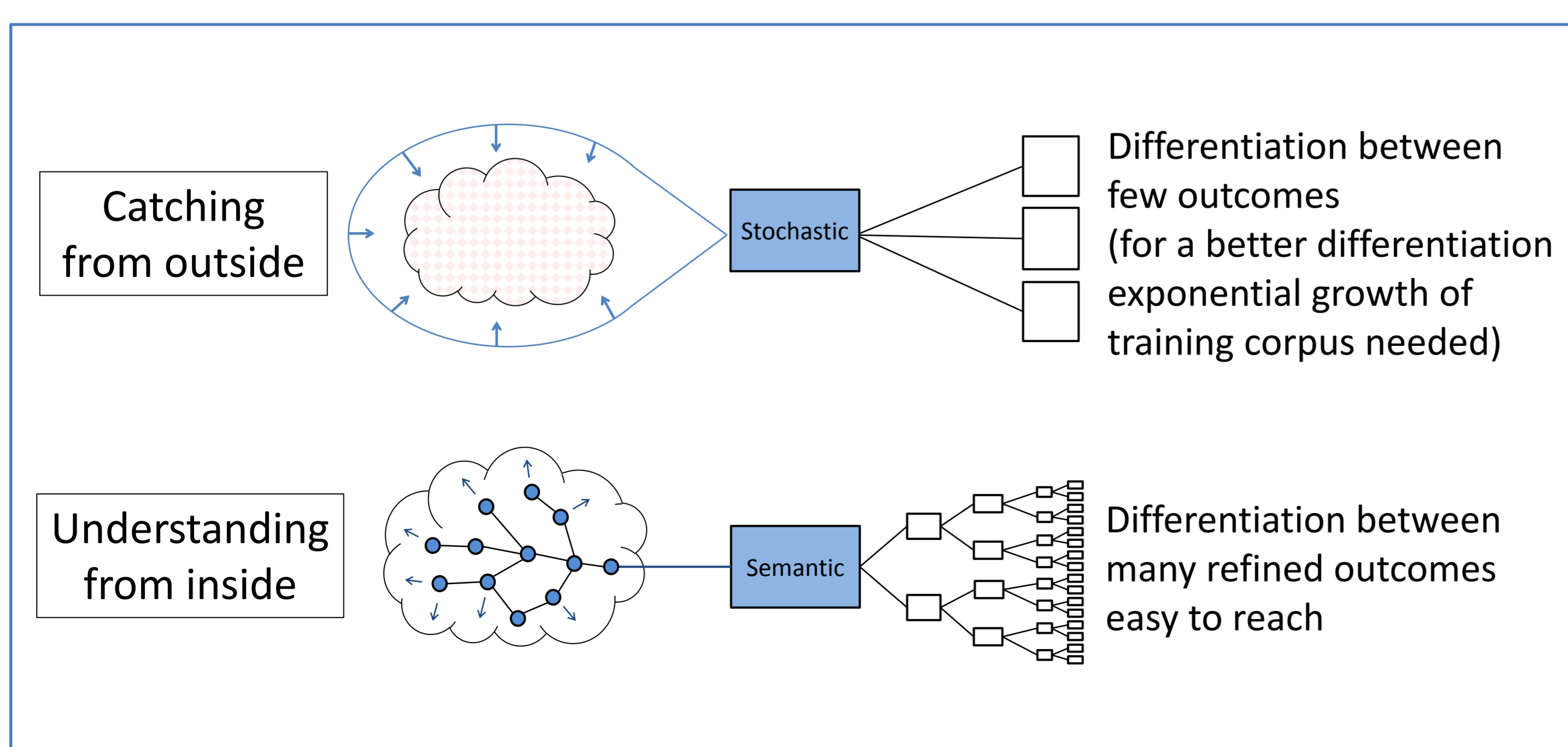
**Concept Molecules** Semfinder technology works with concept molecules (CMs) which capture the semantics. For coding, free text inputs are automatically processed to CMs. CMs are built of atomic concepts, arranged in a structure which represents the relations between the atomic concepts. The actual coding is derived from the information in the molecule.



| Noun phrase ≠ diagnosis phrase (solved by Concept Molecules) |   |
|--|---|
| Linguistics (Input phrase)                                   | Semantics (Interpretation)  |
| Adeno-CA, Colon  | <ul style="list-style-type: none"> <li>diagnosis neoplasia carcinoma adenocarcinoma (1 molecule)</li> <li>malignant (1 diagnosis)</li> <li>ICD-10 c18.9 (1 ICD-10 code)</li> <li>loc intestine large intestine colon</li> </ul>               |
| Adeno-CA, Tinnitus   | <ul style="list-style-type: none"> <li>diagnosis neoplasia carcinoma adenocarcinoma (2 molecules)</li> <li>malignant (2 diagnoses)</li> <li>ICD-10 c80.9 (2 ICD-10 codes)</li> <li>loc tinnitus</li> <li>ICD-10 h93.1</li> <li>loc</li> </ul> |
| Linguistically no difference                                 | Semantically a clear difference   |

## Stochastics and Semantics

Semantic and stochastic methods are complementary in nature.



|  | Stochastics      | Semantics   |
|--|------------------|---|
| <b>Learning Phase</b>                                    | Long (much data) | Long (much expert work)   |
| <b>Recall</b>  | High             | Very high   |
| <b>Noise</b>   | Robust           | Sensitive (LDE, Vocabulary)   |
| <b>Precision</b>   | Medium to High   | Very high   |
| <b>Outcomes</b>  | Few, simple      | Detailed, rich  |
| <b>Multilinguality</b>                                   | Expensive        | Easy access   |
| <b>Process transparency for maintenance /fine tuning</b> | not transparent  | transparent   |
| <b>Further Processing (apart from coding)</b>            | Needs prior work | Ready:<br>→ Semantic Data Repository<br>→ Alerts, Proposals in Clinics<br>→ Clinical Epidemiology |