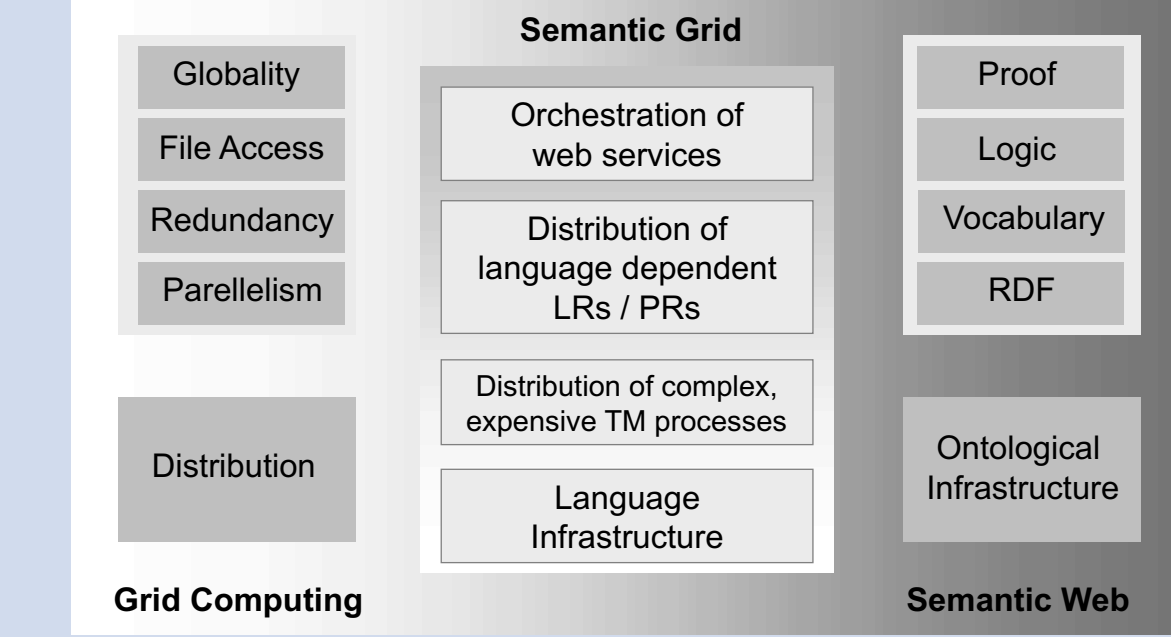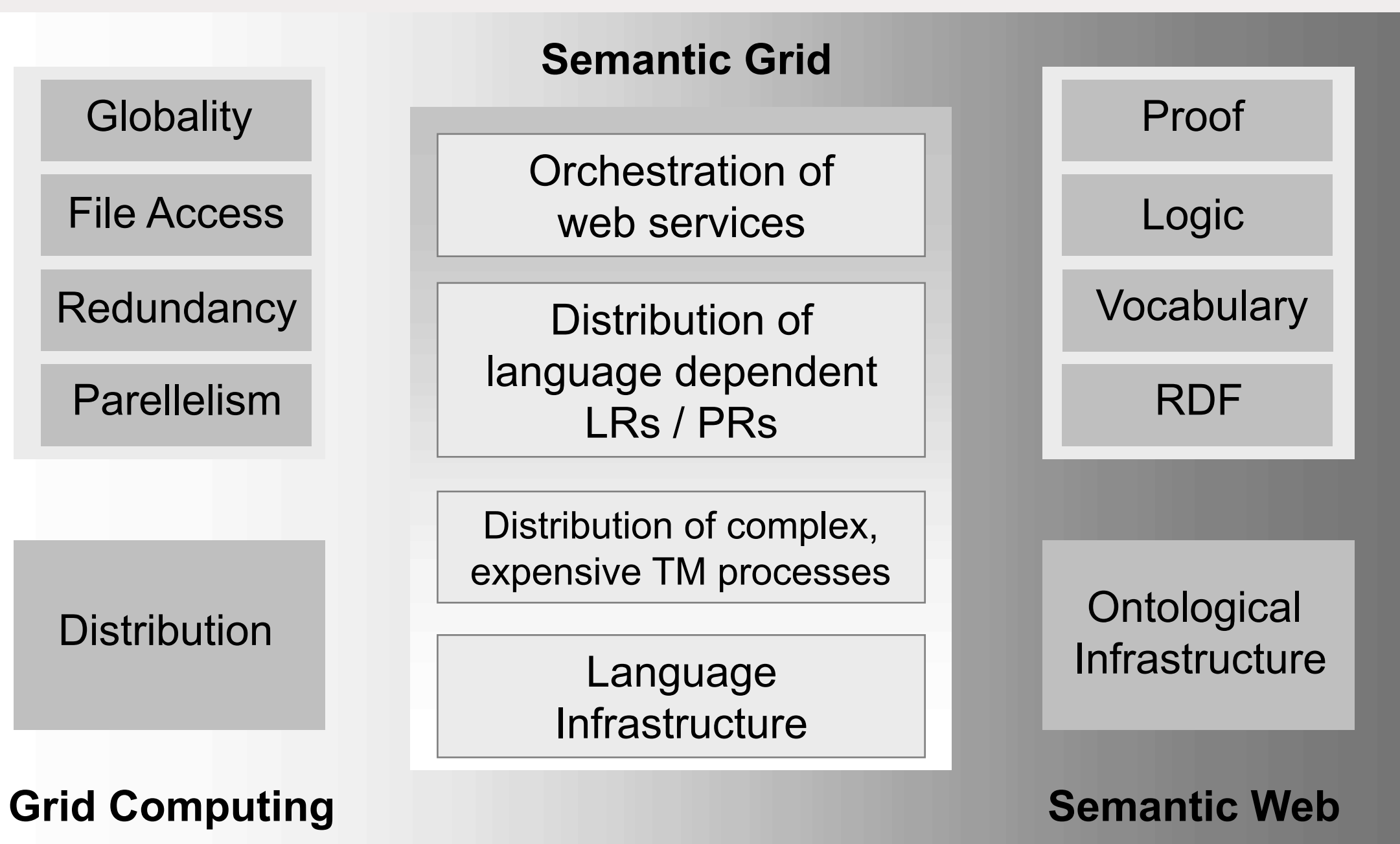# Embedded Distributed Text Mining and Semantic Web Technology

## Daniel Sonntag

German Research Center for Artificial Intelligence, Saarbrücken/Germany, daniel.sonntag@dfki.de
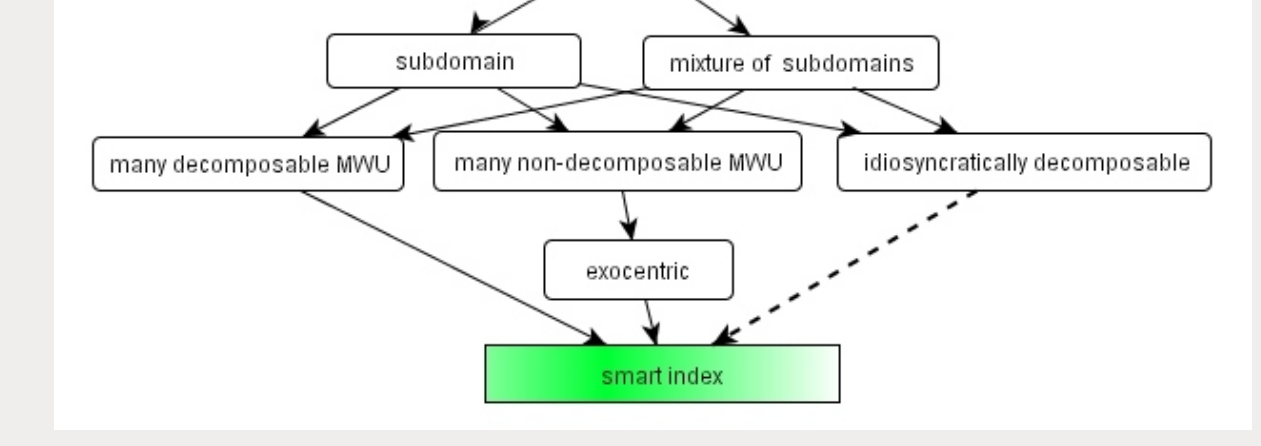
## Common View



## NLP for the Semantic Web

- **Language Infrastructure (1):** Ontology Conceptualisation and Mapping by multi word unit identification. Use frequency and information based methods. Collocation Finding: knock (at) the door, make up, prime minister, etc.
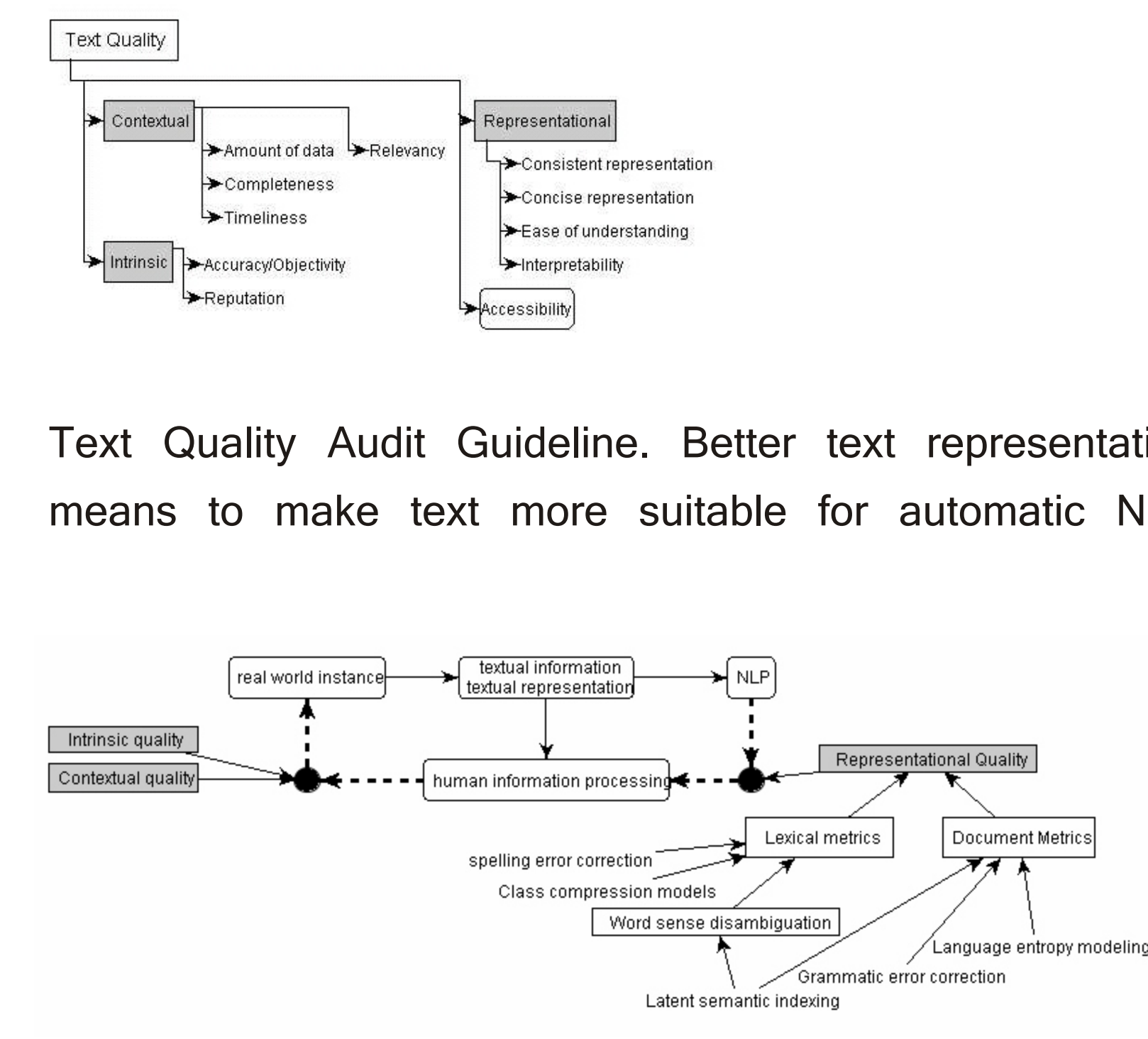


- **Language Infrastructure (2):** Ontology population, (semi-)automatic knowledge markup through Information Extraction, i.e. named entity class tagging and semantic relation tagging.

Ontology Learning provides the instrument for adapting to different application domains.
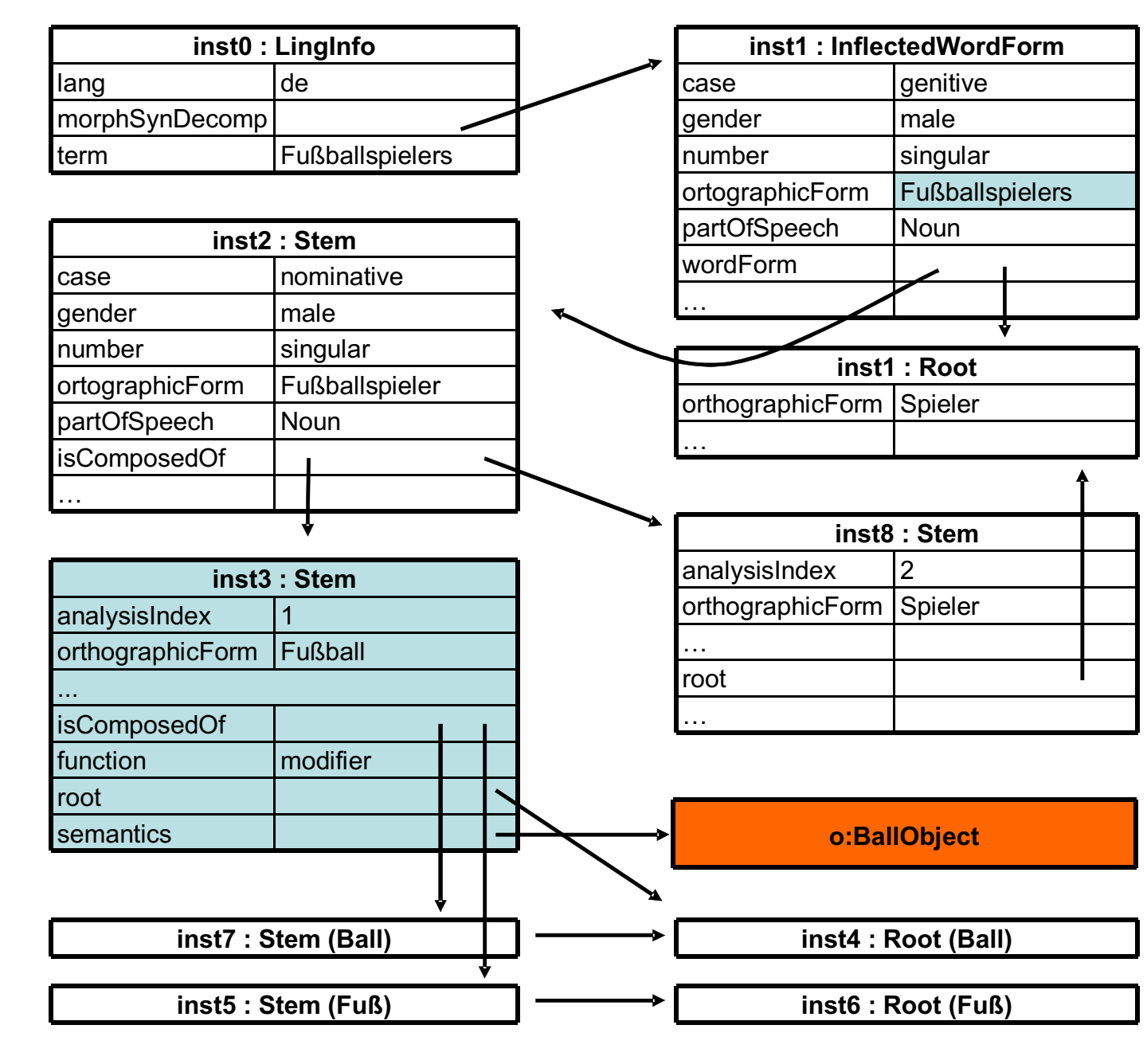
## Representational Data Quality

- Representational data quality includes aspects of data formats in form of concise and consistent representation, and meaning in terms of interpretability and ease of understanding.



- Text Quality Audit Guideline. Better text representation means to make text more suitable for automatic NLP.
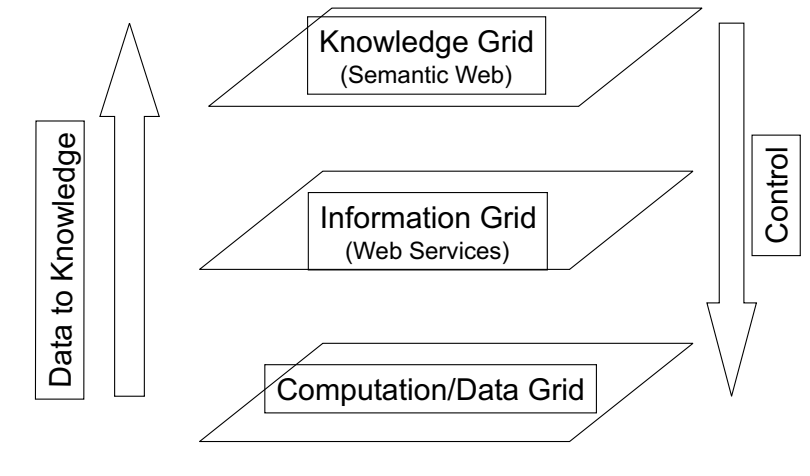


## Semantic Web for NLP

- **Language infrastructure (2):** Semantic Web data structures organise texts in different languages for better representational data quality.

- LingInfo morphosyntactic decomposition instance example of german compound: Fussballspielers (of the football player) according to a domain football player ontology (SWintO).



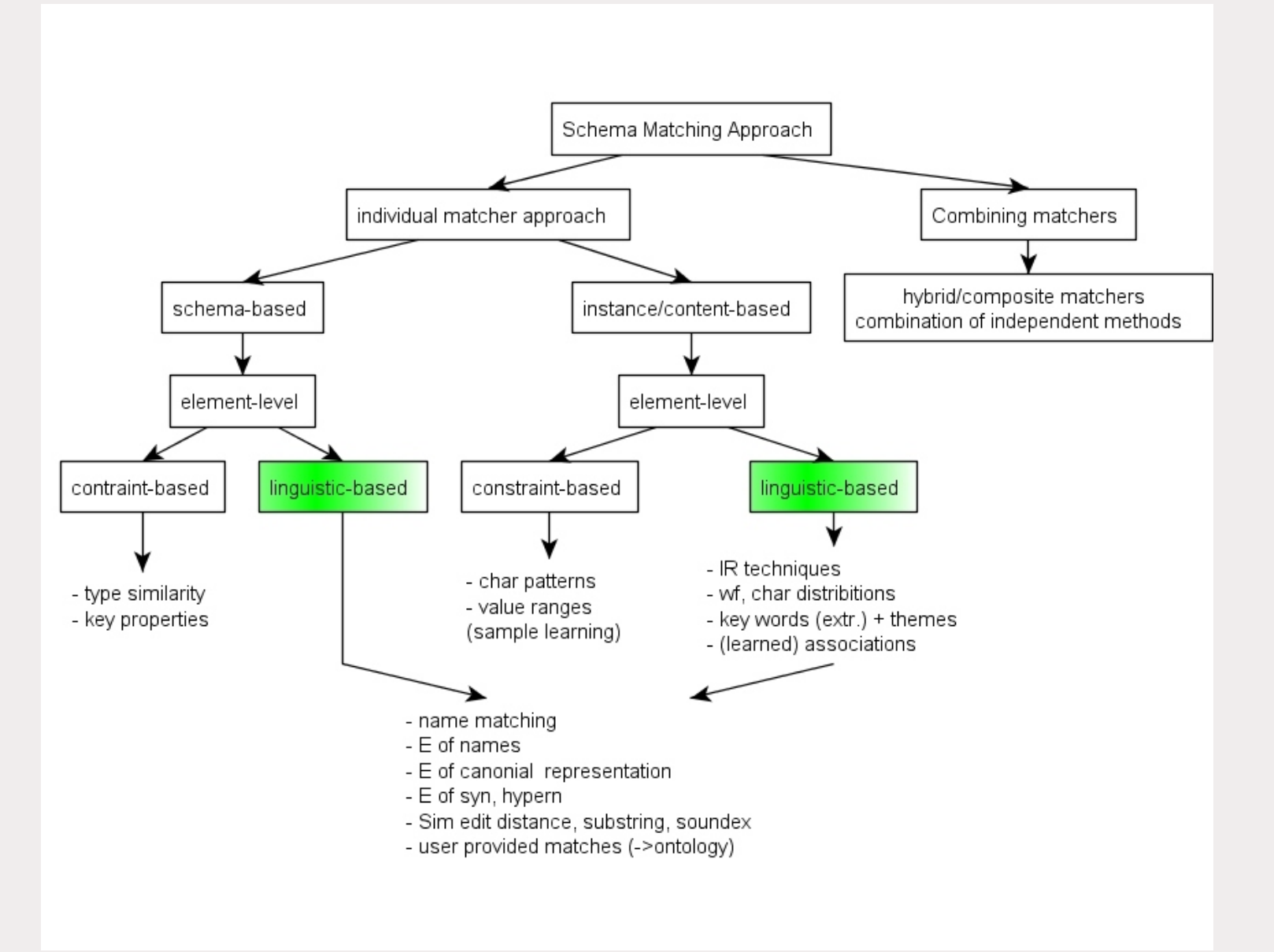- Close Semantic Gap by ontological text annotations.

## Semantic Grid Computing

- Semantic Grid is defined to emphasise the use of Semantic Web technology in the Grid. The Grid is the seamless access to distributed computing and information resources (semanticgrid.org).

- Distribution of Language and Processing Resources over geographically dispersed systems. LRs refer to data resources such as lexica, corpora, thesauri, and ontologies. PRs refer to programmatic or algorithmic sources such as distributed classifiers, POS taggers, NE recognisers, or grammatical parsers.



- Computational Grid vs. Data Grid and three layer Grid abstraction.

- Resource Reasoning: discovery, selection, composition

## Embedded Text Mining Workflow

- (1) Build local (semantic web) services for e.g. text classification by building special purpose classifiers for different natural languages realised through distributed access to a multitude of such classifiers.

- (2) Model PRs and LRs communalities (e.g. equivalence classes) by ontologies.

- (3) Annotate texts and spin your Semantic Web.

- (4) Enable e.g. distributed document access and retrieval to confidential data.

- (5) Decide on different NLP pipelines.

- (6) Deploy embedded knowledge assisted Text Processing services.

- (7) Automatise previous stages by Text Mining, i.e. automatic document markup, automatic ontology population, ontology concept and relation learning, and automatic ontology mapping, i.e. automatic schema mapping.

## Database Technology

- Schema Matching Problems: unknown synonyms/hyponyms, foreign-language data material, cryptic attributes

- Automatic (Linguistic) Schema Matching Approaches



- Use DB operators and ontologies for data scalability and data representation, respectively.

- Expand queries by ontological content operators. These operators are useful in keyword searching, collocation searches, and pattern matching searches.



## Standards

- RDF(S), OWL, SOAP, WSDL, DAMSL, Topic Maps

- MPEG-7 for describing multimedia content data.

- SmartWeb SWEMMA: W3C EMMA Extension for representing ontological queries, results, and status objects.

- Open Grid Services Architecture (OGSA), Service-Oriented Architecture (SOA), ebXML (ISO/TS15000)

- Predictive Model Markup Language (PMML)