



LUND
UNIVERSITY

Question Answering and the development of the *Hajen* System

Pierre Nugues

Joint work with Marcus Klang, Rebecka Weegar, Peter Exner, Juri Pyykkö



Background: IBM Watson

- IBM Watson: A system that can answer questions better than any human
- Video: https://www.youtube.com/watch?v=WFR3IOm_xhE
- IBM Watson builds on the extraction of knowledge from masses of texts: Wikipedia, archive of the New York Times, etc.
- Bottom line: Text is the repository of human knowledge



Goals of *Hajen*

- Build a system that can answer questions in a way similar to IBM Watson
- Use semantic knowledge extracted from Swedish sources: Wikipedia, encyclopedias, *Sydsvenskan*, others?
- Beyond Swedish:
 - Implement a generic, open-source, modular, question answering platform
 - Make it easy to integrate multilingual components
 - Use of full-text indexing, graph-based document models, map-reduce Hadoop ecosystem, dependency injection techniques

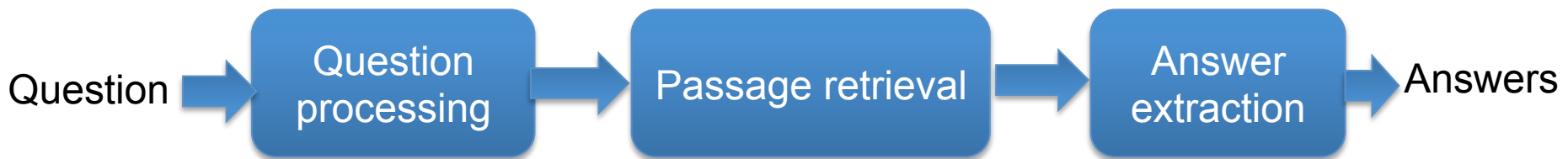


Hajen?

- Part of Vetenskapsrådets frame funding **Det digitaliserade samhället**
- Semantic processing project: *Mot kunskapsbaserad storskalig kunskapsutvinning ur svenskt text*
- Cooperation with Gothenburg University and Chalmers
- The name is a tribute to Ulf Hannerz, the legendary winner of *Kvitt eller dubbelt* in 1957.
- <http://www.svtplay.se/klipp/297574/10-000-kronorsfragan-kvitt-eller-dubbelt>



Question-Answering Architecture (simplified from IBM Watson)



Question parsing and classification:
Syntactic parsing, entity recognition, answer classification

Document retrieval.
Extraction and ranking of passages:
Indexing, vector space model.

Extraction and ranking of answers:
Answer parsing, entity recognition



Corpus: *Kvitt eller dubbelt (KED)*

- In language processing, everything starts with a corpus
- *SVT-klassikern*, a simplified *Jeopardy!* in Sweden



Kvitt eller dubbelt: The Questions

På resa i rymden	Handboll	Aristocats	* Hunden
<p>250 Vad kallar vi rymdfarare med ett annat namn? Astronauter (i Ryssland kosmonauter)</p> <p>500 Vilket land är det hittills enda som skickat upp rymdfärjor? USA</p> <p>1000 Vad svävar i rymden och skickar ner tv-program till parabolantennen? Satelliter</p> <p>2000 Är det sant att man kan se kinesiska muren med blotta ögat från månen? Nej</p> <p>5000 Vilket år landade människan för första gången på månen - 1909, 1969 eller 1999? 1969</p> <p>10000 På vilken planet har vi landat bilar, och hoppas snart få dit människor? Mars</p>	<p>250 Är en handboll större eller mindre än en fotboll? Mindre</p> <p>500 Hur många utespelare per lag finns det på plan? Sex</p> <p>1000 Hur många steg får en handbollsspelare springa utan att studsa bollen? Tre</p> <p>2000 Vad kan man också kalla ett sjumeterskast? Straffkast</p> <p>3000 Vilket svenskt herrlag har vunnit flest SM på senare år - Lugi, Redbergslid eller Skövde? Redbergslid</p> <p>10000 Vem har under många år varit framgångsrik handbollstränare för Sverige? Bengt (Bengan) Johansson</p>	<p>250 Vilken sorts djur är huvudpersoner i filmen Aristocats? Katter</p> <p>500 Vilken fågel har kattungen Marie? Hon är vit</p> <p>1000 Hankatten Thomas hjälper Duchess och hennes ungar. Vad heter han i efternamn? O'Malley</p> <p>2000 Vad är systerna Abigail och Amelia och deras Onkel Waldo för sorts fåglar? De är gäss</p> <p>3000 Vad heter den elaka betjänten som vill ha bort katterna? Edgar</p> <p>10000 Två hundar jagar betjänten. En heter Lafayette. Vad heter den andra, som bestämmer? Napoleon</p>	<p>250 Vad kallas hundens ungar? Valpar</p> <p>500 En arg hund morrar, men vad gör en hund som varnar för fara? Skäller</p> <p>1000 Med vilken kroppsdel visar hunden vanligen att den är glad? Svansen (den viftar på svansen)</p> <p>2000 Med vilket rovdjur tror man hunden är släkt? Vargen</p> <p>3000 Hur ska nosen kännas på en frisk hund? Den ska vara fuktig och kall</p> <p>10000 Vad kallas en hundgård för hunduppfödning med ett annat ord? Kenel</p>

TEKNIK
OCH
UPPPINNINGAR

SPORT
OCH
FRITID

BÖCKER
OCH
FILM

DJUR
OCH
NATUR



Representing the Questions

- We transcribed the 2300 questions and we represented them as an RDF graph.

- A question from the category: “Får jag lov”

Q: Vad heter den argentinska dans som Petra Nielsen dansade i melodifestivalen 2004? A: Tango

- Represented as:

[kvitt:line 1; kvitt:value 250; kvitt:text "Vad heter den argentinska dans som Petra Nielsen dansade i melodifestivalen 2004?"; kvitt:answer "Tango"]

- Using SPARQL, we can extract easily data from the graph and carry out tests.



Passage Retrieval: Textual Resources

- Question answering needs a knowledge source in the form of a collection of documents.
- Wikipedia has a large and growing coverage of topics.
- It is easy to download from dumps.wikimedia.org.
- Is the Wikipedia suitable to answer *Kvitt eller dubbelt* questions?



Passage Retrieval: Indexing

- We segmented the Wikipedia articles into paragraphs: the passages.
- We indexed the passages using Apache Lucene
- Given a question, like:

Vad heter den argentinska dans som Petra Nielsen dansade i melodifestivalen 2004?

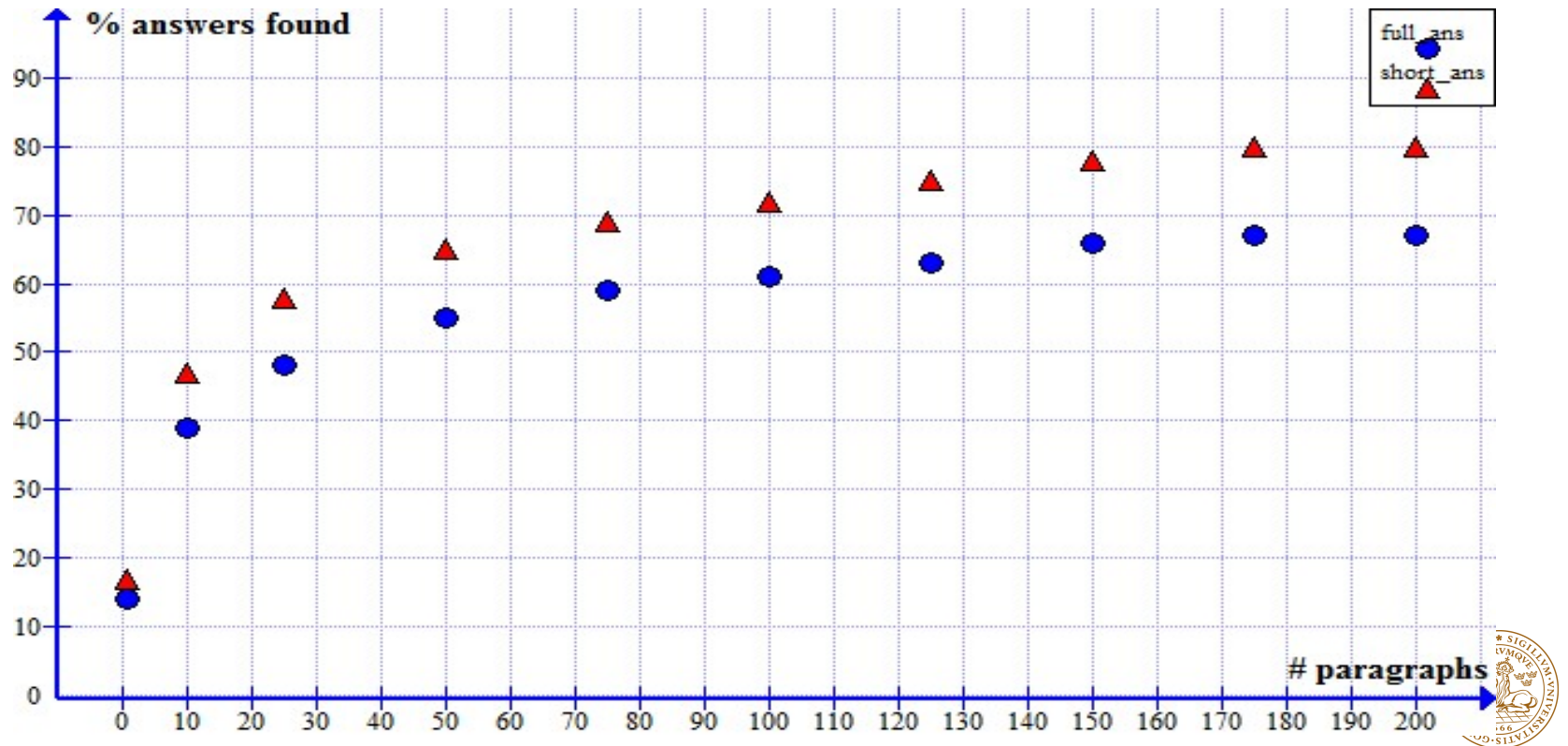
We return the most relevant paragraphs scored using the Lucene's built-in TF.IDF measure.

- [Demonstration: <http://semantica.cs.lth.se:8888/sv/query/passage/>]
- How good is it?



Passage Retrieval: Some Results

We submitted questions from the corpus and we checked if the returned paragraphs contained the answer string.



Joint work with Juri Pyykkö and Rebecka Weegar

Answer Extraction

- Most KED answers are concepts or entities (real or imaginary unique things).
- Corresponds to common nouns or proper nouns in the passages.
- Extractible using a part-of-speech tagger.
- Proper nouns can be classified into categories such as persons, cities, countries, organizations, etc.
- Once extracted, candidates are ranked by increasing frequency.
- This is our baseline answer extraction



Part-of-Speech Tagging

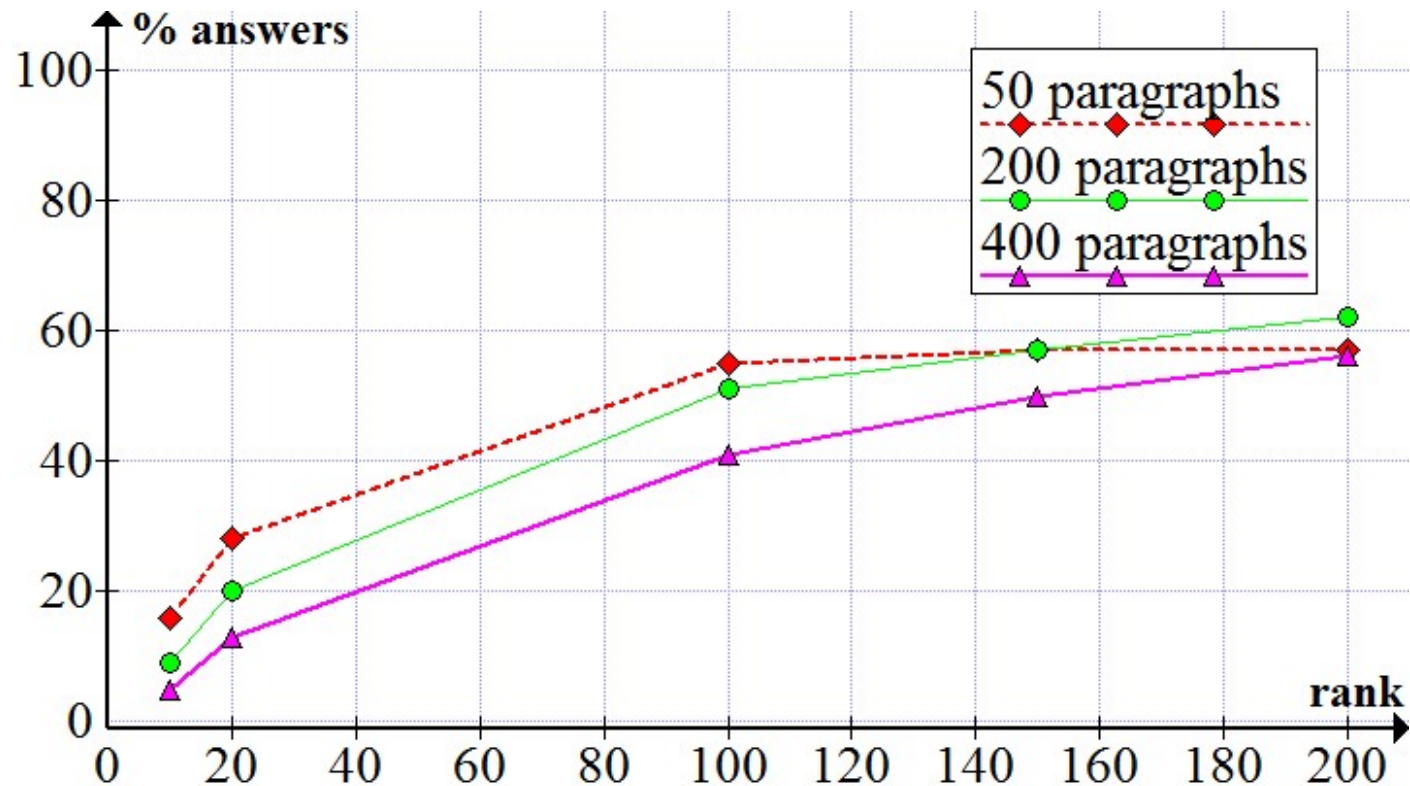
1	Petra	Petra	PM	(person
2	Magdalena	Magdalena	PM	person
3	Nielsen	Nielsen	PM	person)
4	,	,	MID	—
5	ursprungligen	ursprungligen	AB	—
6	Inczèdy-Gombos	inczèdy-gombos	PM	—
7	,	,	MID	—
8	född	föda	PC	—
9		1	1RG	—
10	februari	februari	NN	—
11	1965	1965	RG	—
12	i	i	PP	—
13	Stockholm	Stockholm	PM	(place)
14	.	.	MAD	—



Part-of-Speech Tagging

1	Petra	Petra	PM	(person
2	Magdalena	Magdalena	PM	person
3	Nielsen	Nielsen	PM	person)
4	,	,	MID	—
5	ursprungligen	ursprungligen	AB	—
6	Inczèdy-Gombos	inczèdy-gombos	PM	—
7	,	,	MID	—
8	född	föda	PC	—
9		1	1RG	—
10	februari	februari	NN	—
11	1965	1965	RG	—
12	i	i	PP	—
13	Stockholm	Stockholm	PM	(place)
14	.	.	MAD	—

Answer Extraction: Preliminary Results



Results limited to answers classified as entities.

Joint work with Juri Pyykkö and Rebecka Weegar.



The Mark-1 Question-Answering Engine

Question Answering Demo

Vad heter Sveriges drottning?

Limit (100) ▾

Engine (mark1-full) ▾

Query

1. human
(0.9691119575809469)
2. location
(0.01206060782940241)
3. entity
(0.011618655797206709)
4. numeric
(0.0023526649965215285)
5. description
(0.0013436455495327417)
6. misc
(0.001252322513930337)
7. action
(0.001142055894251944)
8. abbrev
(0.0011180898382073701)

1. drottning Silvia

Score 0.9691119575809469
Alternate answers drottning Silvia (5), Drottning Silvias (3), Drottning Silvia (2), Silvia (2), drottning S (2)
Properties stagger.ne:person (13), wikidata:Q152308 (13)

2. Cardellgatan

Score 0.9691119575809469
Alternate answers Cardellgatan (2)
Properties stagger.ne:place (2)

3. drottning Helena

Score 0.9691119575809469
Alternate answers drottning Helena (2)
Properties stagger.ne:person (2)

4. Drottning Kristina

Demonstration: <http://semantica.cs.lth.se:8888/sv/query/qa>



Answer Type Classification

- From the question text, we can often guess the type of the answer and discard unlikely candidates.
- *Q: Vilket är det enda nordiska land man inte kan köra bil från Sverige? A: Island*
- We annotated a part of the answers from KED and we trained an answer type classifier.
- We use logistic regression and we get probabilities for the answer types.
- Not completely integrated yet.

1. location
(0.9751840864733327)
2. entity
(0.020735020889680874)
3. numeric
(0.0015616698708187907)
4. human
(6.548017611252503E-4)
5. misc
(5.00469009714687E-4)
6. action
(5.004389990210438E-4)
7. abbrev
(4.964844691307781E-4)
8. description
(3.6702852717596216E-4)



Reranking

- A reranker is a module that takes a list of ranked candidates and reorder them using “global” features such the answer type.

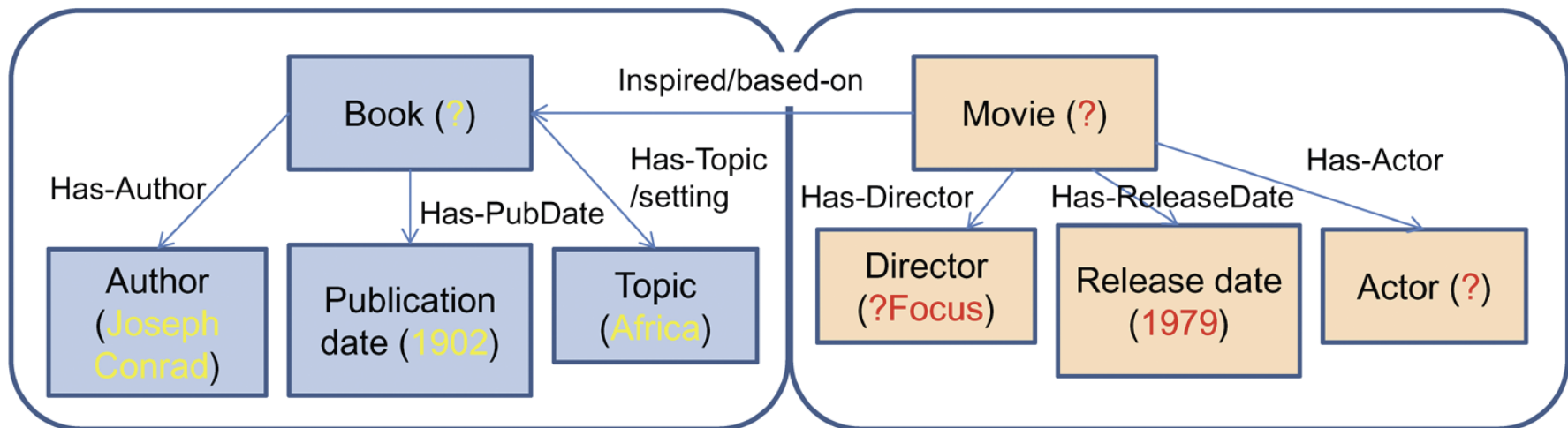
Q: Vilket är det enda nordiska land man inte kan köra bil från Sverige?

Rank	Answer	F/T
1	Europa	false
2	Patrik Budda Andersson	false
3	Sannas	false
...		
67	Island	true

- Using the corpus, it is possible to train a binary classifier with features such as the predicted answer type, word type
- Improves considerably the results (Not yet in the Mark-X line)

Structured Data and Inference (From IBM Watson)

WAR MOVIES: A 1902 Joseph Conrad work set in Africa inspired this director to create a controversial 1979 war film.



Named Entity Processing

- A first step to inference is to identify and disambiguate named entities.
- Two techniques: named entity recognition (NER) and named entity disambiguation (NED)

Göran Persson var statsminister mellan åren 1996 och 2006 samt var partiledare för Socialdemokraterna

- Google's motto *Things not Strings*



Named Entity Recognition

- Named entity recognition (NER) tags the words that correspond to named entities, possibly with their type.
- Common types include person, organization and location.
 - **Göran Persson** is a person.
 - **Socialdemokraterna** is an organization.
- Techniques to carry out named entity recognition use extensions of part-of-speech tagging



Named Entities are Ambiguous

Göran Persson kan syfta på

- Göran Persson (född 1949), socialdemokratisk partiledare och svensk statsminister 1996–2006
- Göran Persson (född 1960), socialdemokratisk politiker från Skåne
- Göran Persson (musiker), svensk proggmusiker
- Jöran Persson, svensk ämbetsman på 1500-talet



Named Entity Disambiguation

Named entity disambiguation (NED) links words to entities: *strings to things*.

Göran Persson var statsminister mellan åren 1996 och 2006 samt var partiledare för **Socialdemokraterna**.



Wikipedia and Wikidata

- Each page in Wikipedia (concepts and entities) has a unique identification number across all the languages, a sort of *personnummer*.
- The Q-number repository (wikidata) is a growing source of structured data.
 - Göran Persson (statsminister): Q53747
 - Göran Persson (riksdagsledamot): Q5626648



Wikidata Content

In other languages

français	Göran Persson
	homme politique suédois
svenska	Göran Persson
	Sveriges statsminister 1996-2006, Sveriges finansminister 1994-1996 samt Sveriges skolminister 1986-1991
suomi	Göran Persson
	<input type="text" value="enter a description in suomi"/>

Statements

sex or gender	<input type="checkbox"/> male
	▶ 3 references
LCNAF identifier	<input type="checkbox"/> n98083850 
	▶ 1 reference

In other languages

français	<input type="text" value="enter a label in français"/>
	<input type="text" value="enter a description in français"/>
svenska	Göran Persson (född 1960)
	<input type="text" value="enter a description in svenska"/>
suomi	<input type="text" value="enter a label in suomi"/>
	<input type="text" value="enter a description in suomi"/>

Statements

sex or gender	<input type="checkbox"/> male
	▶ 1 reference
instance of	<input type="checkbox"/> human
	▶ 1 reference



NEDforia: A Disambiguation Tool

Named Entity Disambiguation Result

Göran Persson var statsminister mellan åren 1996 och 2006 samt var partiledare för Socialdemokraterna.

Result information

Named Entity Disambiguation Candidates

Named Entity

Göran Persson

rank	uri	rank-value	commonness
1.	wikidata:Q53747	0.9999942697886333	0.9809523809523809
2.	wikidata:Q6042900	0.753557930607055	0.0126984126984127
3.	wikidata:Q6043257	0.7290883542424511	0.0031746031746032
4.	wikidata:Q2625684	0.7246375610550542	0.0031746031746032
5.	wikidata:Q5626648	0.6595001745138378	0

Nedforia Marcus Klang 2014

Demonstration: <http://semantica.cs.lth.se:8888/sv/ned/query>



LUND
UNIVERSITY

Mark 1 Architecture

Webforia (HTTP frontend)

Toolforia (CLI)

Nlpforia

Language
(Dependency Injector)

Parsers – Wikimarkup, ConLL

Disambiguation – Detectors and disambiguators

Data – Fulltext, Passage search, Key/Value storage

Classification – Short text (based on LibLinear)

QA – Rankers, Extractors, Hypothesizers, Combiners

NLP Tools – Tokenizers, NER, Segmenters, POS Taggers

Global Dependency Injection

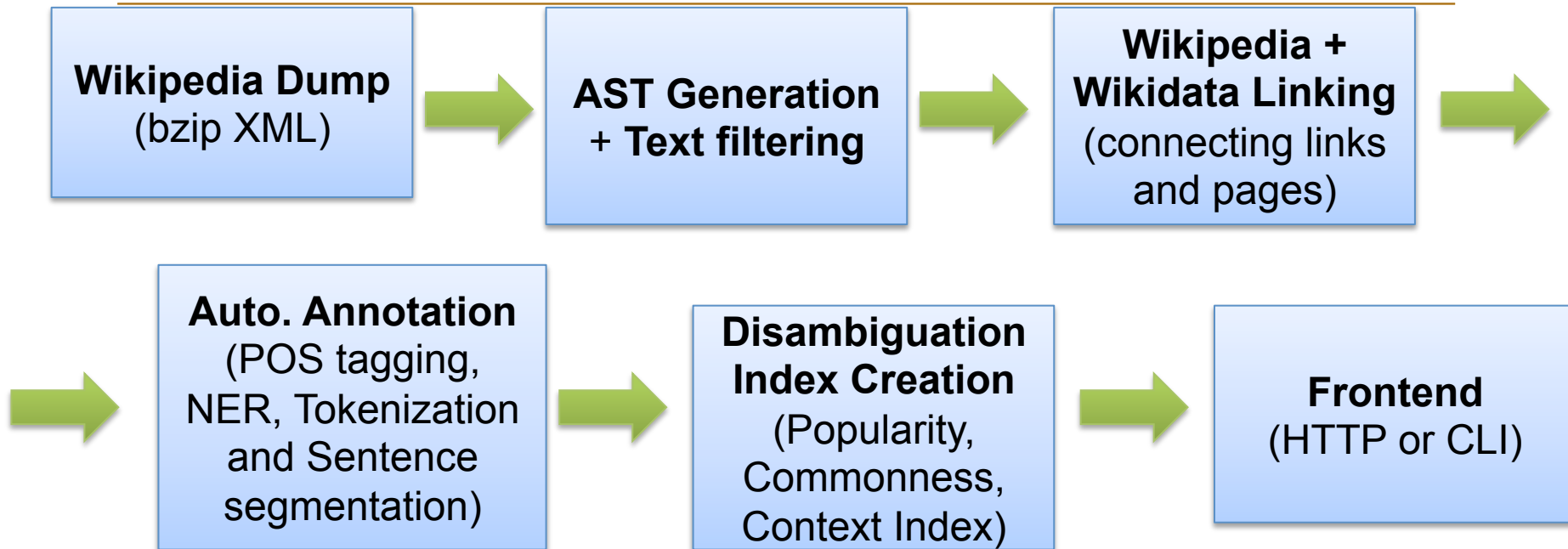
Document Model (Graph based)

Mark 1: Features

- **Document model**
 - Directed property graph
 - Built-in small query engine
 - Extends or re-uses ideas from KOSHIK (from Peter Exner)
- **Multilingual and modular**
 - **Idea:** A Language is nothing more than an assemblage of language specific components with models.
 - Constructor based dependency injection (Guice)
- **Designed for prototyping and experimentation**
- **Open source**



NEDforia: The NED Pipeline



Auto annotation was based on Stagger and Maltparser with respective models trained on the Stockholm-Umeå Corpus (SUC)



NEDforia: Disambiguation Results

- Small test corpus consisting of 10 news articles from mixed sources:
Aftonbladet, Svd, DN, Nyteknik, HD, IDG and Dagens Industri
- We extracted and annotated 201 named entities, preliminary results.

Recall	Precision	F1
0.696	0.748	0.721



Cross-Language Extraction

English frame		
SBJ	Einstein	A0
VERB	received	receive.01
NMOD	the	
NMOD	1921	
OBJ	Nobel Prize	A1
TMP	in	
PMOD	Physics	

Swedish frame		
SS	Einstein	A0
VERB	fick	få.01
DT	1921	
DT	års	
OO	Nobelpris	A1
ET	i	
HD	fysik	

Koshik

- Koshik: End-to-end framework to process multilingual documents
- Based on Hadoop and our Crafoord cluster
- Koshik is a talking elephant (Korean)



Koshik Architecture

- Koshik uses the Avro binary format to serialize the documents.
- Avro is designed for Hadoop and allows other data warehousing tools to directly query the documents.
- Koshik can be queried directly through Hive, which offers an SQL-like query language called HiveQL.
- Available at <https://github.com/peterexner/KOSHIK>



Querying with HIVE

Number of articles:

```
> SELECT count(identifier) from koshibdocs;
```

```
Job 0: Map: 920 Reduce: 1 Cumulative CPU: 35243.42 sec HDFS Read: 248123756147
```

```
HDFS Write: 8 SUCCESS
```

```
OK
```

```
4012291
```

Number of tokens:

```
> SELECT count(ann) FROM koshibdocs LATERAL VIEW explode(annotations.layer)
```

```
annTable as ann WHERE ann LIKE '%Token';
```

```
Job 0: Map: 920 Reduce: 1 Cumulative CPU: 35199.98 sec HDFS Read: 248123756147
```

```
HDFS Write: 11 SUCCESS
```

```
OK
```

```
1485951256
```

Number of nouns:

```
> SELECT count(key) FROM (SELECT explode(ann) AS (key,value) FROM (SELECT ann
```

```
FROM koshibdocs LATERAL VIEW explode(annotations.features) annTable as ann) annmap)
```

```
decmap WHERE key='PPOS' AND value LIKE 'NN%';
```

```
Job 0: Map: 920 Reduce: 1 Cumulative CPU: 46694.33 sec HDFS Read: 248123756147
```

```
HDFS Write: 10 SUCCESS
```

```
OK
```

```
476161210
```



A Snapshot of the Future

When wireless is perfectly applied, the whole earth will be converted into a huge brain, which in fact it is, all things being particles of a real and rhythmic whole. We shall be able to communicate with one another instantly, irrespective of distance. Not only this, but through television and telephony we shall see and hear one another as perfectly as though we were face to face, despite intervening distances of thousands of miles; and the instruments through which we shall be able to do this will be amazingly simple compared with our present telephone. A man will be able to carry one in his vest pocket.

Nikola Tesla, *Colliers*, January 30, 1926

