



Association for
Computing Machinery



Information to Wisdom: Commonsense Knowledge Extraction and Compilation

Simon Razniewski

Max Planck Institute for
Informatics, Germany



mpi max planck institut
informatik

Niket Tandon

Allen Institute for AI,
Seattle, WA, USA

Ai2 Allen Institute for AI

Aparna S. Varde

Montclair State University,
Montclair, NJ, USA



MONTCLAIR STATE
UNIVERSITY

**Research Tutorial at ACM WSDM conference
(Web Search and Data Mining) - March 2021**

Outline

09:00 IST	15 min	1. Introduction to commonsense knowledge (Simon)
09:15 IST	35 min	2. Text extraction (Simon)
09:50 IST	10 min	<i>Break</i>
10:00 IST	20 min	3. Multimodal knowledge (Niket)
10:20 IST	30 min	4. Deep learning-based techniques (Niket)
10:50 IST	10 min	<i>Break</i>
11:00 IST	25 min	5. Evaluation of the acquired knowledge (Aparna)
11:25 IST	20 min	6. Highlights, outlook and open issues (Aparna)

Outline

1. Introduction

1. What is CSK?
2. Why is it important?
3. How to represent it?
4. What makes it challenging?

What is commonsense knowledge?

Definition 1 (by commonality):

Knowledge shared by most humans

- Possible qualifications
 - Across cultures
 - From early in life (=children)
- E.g., elementary school exam questions
 - <http://data.allenai.org/ai2-science-questions>

What is commonsense knowledge?

Definition 2 (by knowledge type):

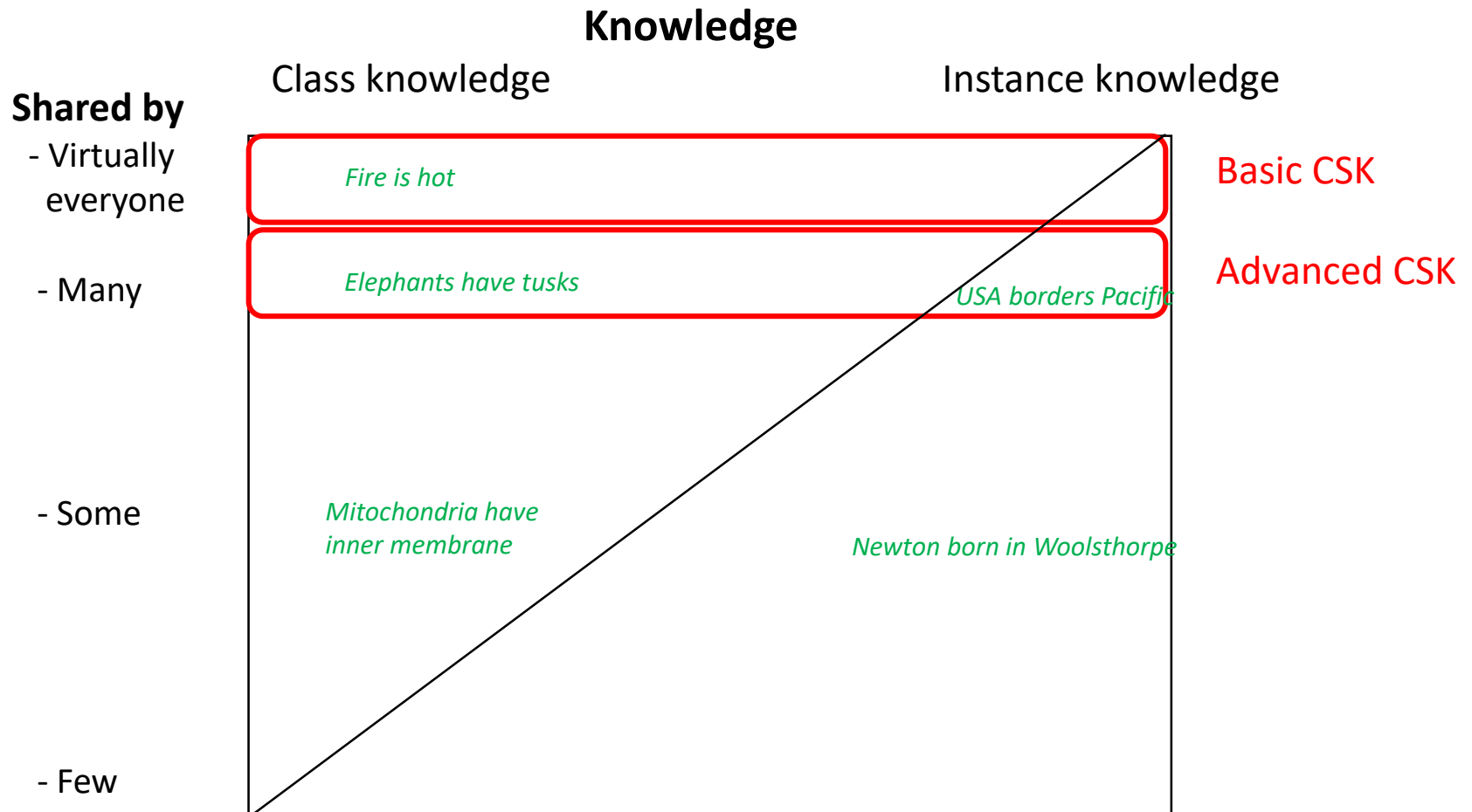
Knowledge about concepts and events

- Concepts: *City, footballer, organization*
- Events: *Football match, birthday party*
- Differentiation from encyclopedic knowledge on instances
 - Instances: *Jerusalem, Ronaldo, Manchester United*

Definition Pro/Con

- Definition 1 (by commonality):
 - *Popsicle, is, frozen – only known in North America*
 - *Lion is dangerous/cute - depends whom you ask*
 - Inclusion/exclusion decision challenging
 - Definition 2 (by knowledge type):
 - *Apple MacBook, Ford Model T*
 - Class/instance not trivial to separate
 - *USA borders Pacific Ocean – excluded as instance knowledge*
 - *Mitochondria, hasPart, inner membrane – not common knowledge*
 - Open-ended
- See part 5 (evaluation) - use ranking

Definition: Merger



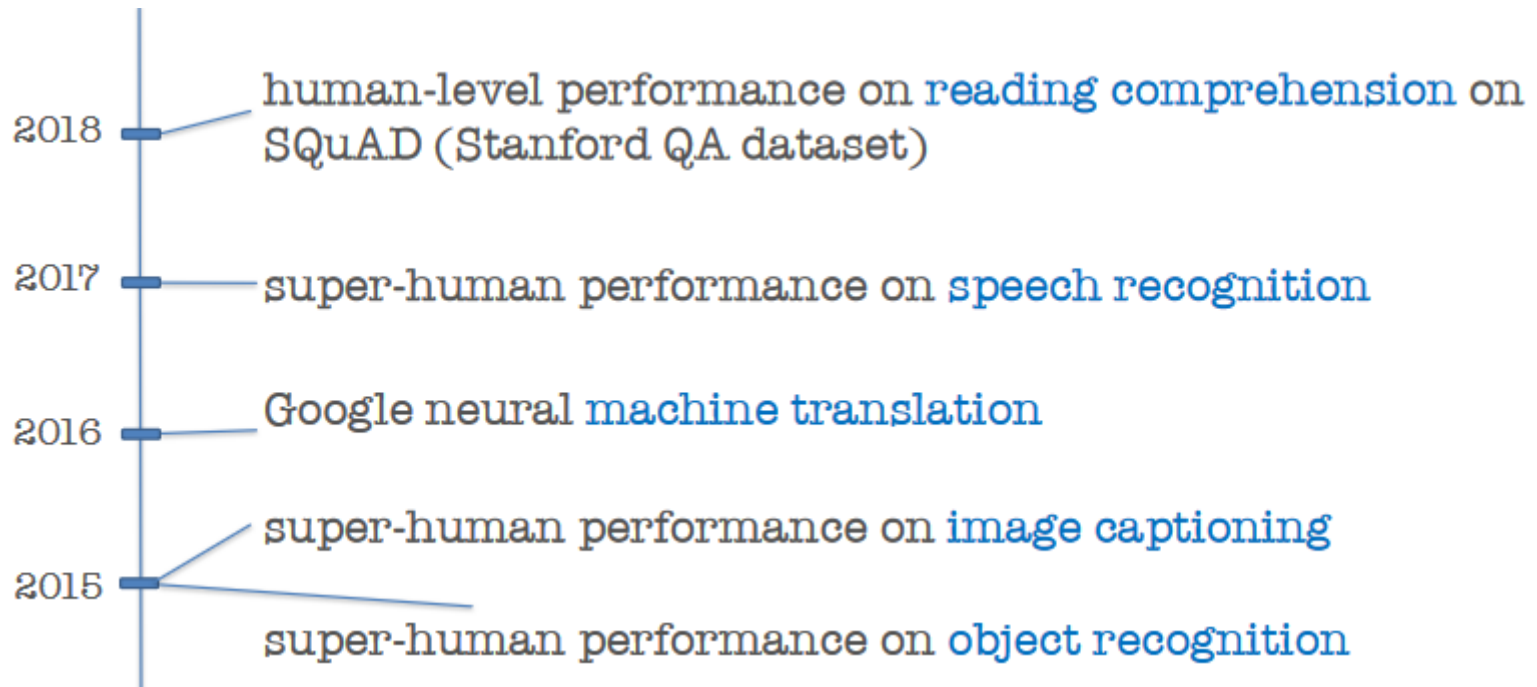
Examples of CSK

- Taxonomical
 - *Elephant, isA, mammal*
- Properties
 - *Elephant, lives in, Savanna*
- Parts
 - *Elephants, hasPart, trunk*
- Measures
 - *Adult elephant, weight, ~2..5 tons*
 - *Elephant, lifespan, ~60 years*
- Activities
 - *Seeing elephant, requires, go to zoo*
 - *Go to zoo, subevent, buy ticket*
 - *Go to zoo, typicalDuration, 2 hours*

What we do not cover

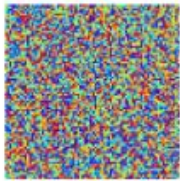
- **Lexical knowledge**
 - Word senses, synonymy, ...
 - WordNet as prime example: <https://wordnet.princeton.edu>
- **Taxonomic knowledge**
 - Good coverage in lexical projects
 - WordNet, WebIsADB, ...
- **Encyclopedic KBs**
 - Wikidata
 - Structured sister project of Wikipedia, mostly focused on instance knowledge
 - Knowledge on concepts slowly growing, though limited set of useful predicates
 - Recent analysis: [Commonsense Knowledge in Wikidata, Ilievski et al., Arxiv, 2020]
 - NELL, DBpedia, YAGO, ...
 - Similar instance focus

Why CSK? Amazing progress without





+



=



Giant panda
Object
Recognition

Gibbon

Szegedy et al,
2014....



VQA

Jabri et al,
2017



A horse standing in the grass.

Captioning

MacLeod
et al, 2017

How are you
doing?



I don't know.

Dialogue

Li et al,
2016



I don't know. I
don't know. I
don't know.

Open-ended

Generation

Holtzman
et al, 2018

.... Nikola Tesla moved to
Prague in 1880. ... **Tadakatsu**
moved to Chicago in 1881.

Where did Tesla move in
1880? **Chicago**

QA

Jia et al,
2017

Solving only a "dataset"
without solving the underlying "task"!

Importance of CSK

Reusable and **scrutable** asset for a range of AI tasks

- **Reusable:**
 - CSK can be plugged into a range of tasks, e.g., QA, dialogue, object recognition, text generation, ...
 - Contrasts with typical end-to-end learning
- **Scrutable:**
 - Humans can inspect, add and **remove** content
 - Relevant in applications where errors are costly
 - Relevant in applications at risk of bias/discrimination
 - Humans can inspect discrete statements used for reasoning
 - Relevant for debugging complex downstream use cases
 - Contrasts with end-to-end learning and pretrained language models

Knowledge representation challenges

- Encyclopedic KBs: Typically binary truth notion
 - *Trump, born in, NY*
 - *House of Cards, producer, Netflix*
 - *New York, mayor, Bloomberg, [2002-2013]*
- CSK: Generalizes across subjects
 - *Lions, have, manes* - percentage?
- Fuzzy time notion
 - *Lions, drink, milk* - when?
- Spatial and cultural context
 - *Lion, is, cute*
 - *Elk, usedFor, transport*

KR - state of the art

- Expressive proposals exist
 - Modal, epistemic, episodic logic
- Instantiation hard
 - Sparse realization in natural language
 - Correct extraction nontrivial
- Most projects:
Pragmatic choice of (subject, predicate, object) triples with a single score

Lion, hunts, zebra – 0.73

Lion, drinks, milk – 0.45

Triples and done?

- Still major design decisions left!

1. Fixed or open set of predicates
2. Subject range
3. Object range

- Fixed vs. open predicates

- E.g., ConceptNet: ~25 predicates (isCapableOf, requires, isA) vs. TupleKB ~1000 textual phrases

- Subjects: Strings or disambiguated terms?

- Lynx vs. lynx vs. lynx

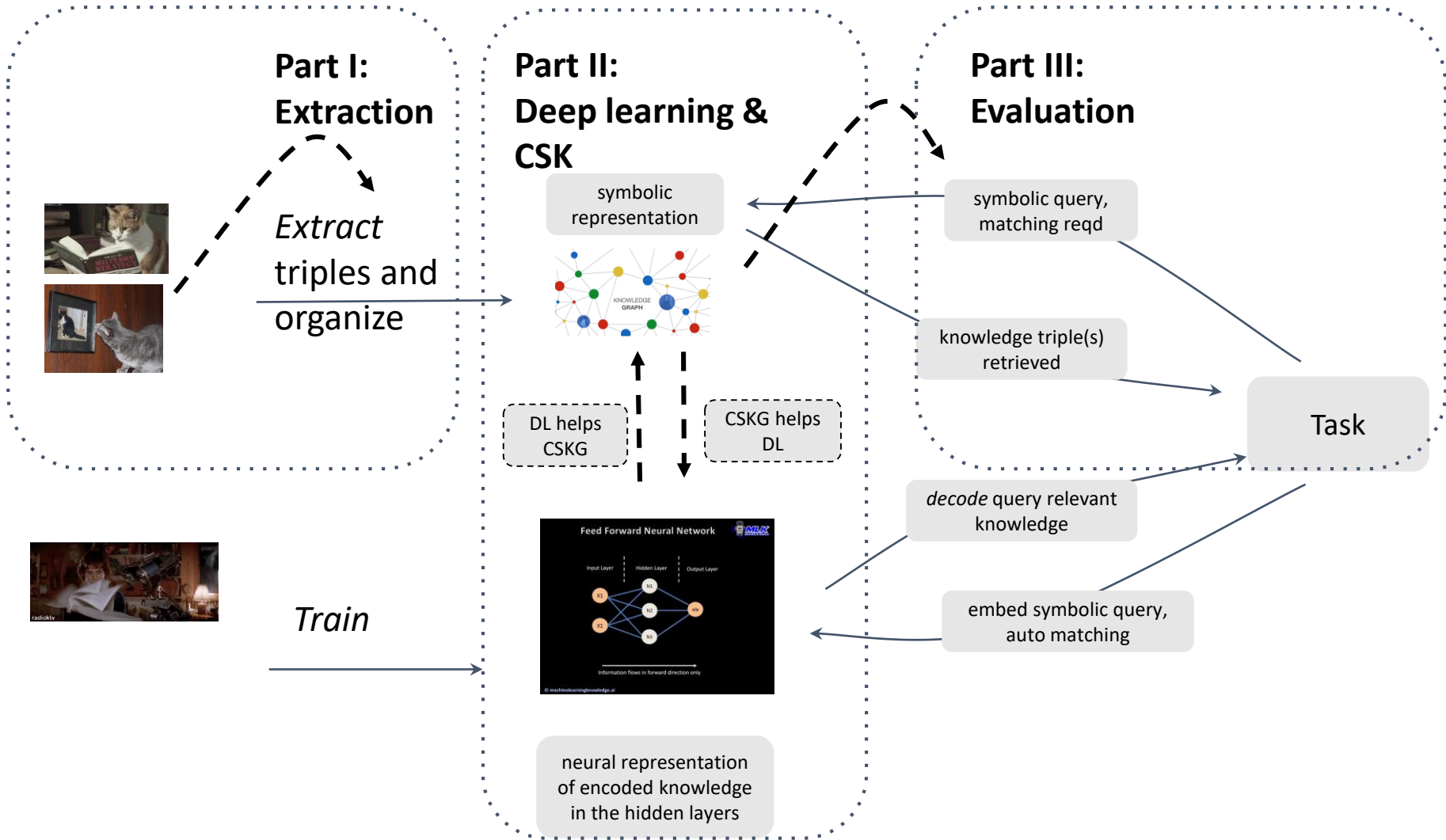


- Granularity and modifiers

- Elephant, Foraging elephant? Newborn elephant?

- Objects: Entities or open phrases?

- Politician, isCapableOf, promise that impossible things will happen



Outline

09:00 IST	15 min	1. Introduction to commonsense knowledge (Simon)
09:15 IST	35 min	2. Text extraction (Simon) - Overview - Recipe - Case studies
09:50 IST	10 min	<i>Break</i>
10:00 IST	20 min	3. Multimodal knowledge (Niket)
10:20 IST	30 min	4. Deep learning-based techniques (Niket)
10:50 IST	10 min	<i>Break</i>
11:00 IST	25 min	5. Evaluation of the acquired knowledge (Aparna)
11:25 IST	20 min	6. Highlights, outlook and open issues (Aparna)

Overview

- **Earliest projects** on CSKB construction were **manually** authored (Cyc, ConceptNet)
- Challenges in scale
 - Atomic: ~100k\$ annotator expenses
- **Automated information extraction** and KB construction field with **long history**
 - Focus traditionally on crisp ``encyclopedic'' knowledge (cf. DBpedia, YAGO, NELL, DeepDive, ...)
- **Can we use automated IE and KBC for CSK?**



Graham Neubig

@gneubig

Following



One commonly cited argument about the difficulty of learning common-sense reasoning is that "no-one writes down common sense". A counter-argument is "well, the web is big": [instructables.com/id/How-To-Open...](https://www.instructables.com/id/How-To-Open...)

How to Open a Door

Step 1: Locate Desired Door

Step 2: Locate Door Handle or Knob

Step 3: Turn Knob or Handle and Pull or Push

Challenges of automated CSKB construction

- **Underspecified text semantics**
 - “Lions attack humans” – all/some/all the time/once/..?
- **Reporting bias**
 - “woman kills” vs. “woman breathes” – 1.5M vs. 0.1M web search results
 - “pink elephant” vs. “grey elephant” – 6.9M vs. 1.9M web search results
- **Sparse observations of quadratic+ space of possible statements**
 - Do computer programmers drink water?
- **Noise and polysemy**
 - Pigs can fly - idiom
 - Lynx: Constellation, web browser, animal

Outline

09:00 IST	15 min	1. Introduction to commonsense knowledge (Simon)
09:15 IST	35 min	2. Text extraction (Simon) - Overview - Recipe - Case studies
09:50 IST	10 min	<i>Break</i>
10:00 IST	20 min	3. Multimodal knowledge (Niket)
10:20 IST	30 min	4. Deep learning-based techniques (Niket)
10:50 IST	10 min	<i>Break</i>
11:00 IST	25 min	5. Evaluation of the acquired knowledge (Aparna)
11:25 IST	20 min	6. Highlights, outlook and open issues (Aparna)

(Textual) information extraction

- Textual information extraction long attention in [KBC/NLP](#)
- Idea: [Exploit patterns/commonalities](#) in natural language in order to extract commonsense knowledge
 - *Lynx eat hares*
 - *Elephants eat grass*
- <s> eats <o> - pattern for (s, diet, o)
- [Generic design points](#)
 1. Sources
 2. Extraction method
 3. Type of contextualization
 4. Consolidation method

Design point 1 – Source choices

- *“Where to extract from?”*

- Wikipedia
- Books and other dedicated sources
 - ARC science corpus
 - Project Gutenberg
- Web search
- Forums
 - Reddit
 - Quora
 - Yahoo Answers
- Search engine query logs
- Web crawls
 - ClueWeb
 - CommonCrawl
- ...



Precision
Coherence

Recall
Redundancy

Extraction source - considerations

- (CS)KB projects stand and fall with source selection
- Precision: *Topic-specific sources >> random web*
 - Event knowledge – Wikihow [HowToKB, WWW 2017]
 - Cultural knowledge – Movie scripts [Knowlywood, CIKM 2015]
 - Science knowledge – Science textbooks [GenericsKB, Arxiv 2020]
- Frequency signals may be stronger from general web dumps, but considerable noise
- Intermediate setting: *Targeted web search* [TupleKB, Ascent]

Design point 2 – Extraction method options

- *“How to extract”*

1. **Manual patterns** [WebChild, WSDM 2014]
 - Hearst patterns etc.
2. **Co-occurrence** [DoQ, ACL 2019]
 - Window, same sentence, ...
3. **Open information extraction** [TupleKB, Quasimodo, Ascent]
 - Any verb phrase
4. Relation-specific **supervised learning**

Extraction method - considerations

- Preferred method depends on desired knowledge representation
 - E.g.,
 - Few non-overlapping relation → Co-occurrence
 - Moderate relations → Supervised extractors
 - Many relations → OpenIE
- Has implications downstream
 - Extraction confidences (supervised extractors) for quantitative contextualization
 - Text context for qualitative contextualization
 - OpenIE with many unspecific extractions

Design point 3 – Contextualization

“What do we annotate statements with?”

1. Observation frequency [WebChild 2.0, DoQ]
 - *Elephant, has, tusks, 155*
 - *Elephant, has, tail, 84*
2. Quantitative [0,1] truth labels [TupleKB, Quasimodo]
 - *Elephant, lives in, group, 0.87*
3. Qualitative truth labels [Ascent]
 - *Elephant, lives in, group, **temp**: during wet season*
 - ***Subgroup**: Female elephant, lives in, group*

Contextualization - considerations

- **Frequencies** trivial to interpret, but do not qualify degree of truth
- **Quantitative truth labels** nontrivial semantics
- **Qualitative labels** easier to interpret, but harder to compare
- **Expressive proposals** from KR exist (e.g., modal logics)
 - Actual implementation not easy
 - Sparse realization in natural language
 - Correct extraction nontrivial

Design point 4 – Consolidation

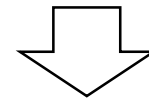
“What do we do with redundant and competing extractions?”

- Similar statements may be seen several times
- Redundancy and contradictions may require additional inference
- Common consolidation methods
 1. Keep all [DoQ]
 2. Frequency cutoff [Ascent]
 - E.g., at least seen 5 times
 3. Per-statement consolidation [TupleKB, Quasimodo]
 - Feature-based classification/ranking
 4. Joint consolidation [WebChild, Dice, Ascent]
 - E.g., BERT-based clustering, MaxSAT, ...

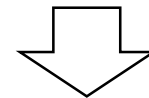
elephant is capable of...

- en carry a trunk →
- en forget to go on the paper →
- en lift logs from the ground →
- en to lift the tree →
- en remember water sources →

Cats, are, solitary
Lions, live in, groups



Lions, are, cats



Consolidation - considerations

- Redundancy challenge and blessing
- Exploiting redundancy requires strong text similarity/entailment modules
- Previous projects often stuck to per-statement consolidation due to lack of strong similarity/entailment modules
- Recent advances on pretrained LMs give hope for joint consolidation (see e.g., Dice, Ascent)

Outline

09:00 IST	15 min	1. Introduction to commonsense knowledge (Simon)
09:15 IST	35 min	2. Text extraction (Simon) <ul style="list-style-type: none">- Overview- Recipe- Case studies
09:50 IST	10 min	<i>Break</i>
10:00 IST	20 min	3. Multimodal knowledge (Niket)
10:20 IST	30 min	4. Deep learning-based techniques (Niket)
10:50 IST	10 min	<i>Break</i>
11:00 IST	25 min	5. Evaluation of the acquired knowledge (Aparna)
11:25 IST	20 min	6. Highlights, outlook and open issues (Aparna)

Representative projects

1. **Webchild 1.0 [Tandon et al., WSDM 2014]**
 - Disambiguated noun-adjective pairs
2. Quasimodo [Romero et al., CIKM 2019]
 - Salient general triples
3. DoQ [Elazar et al., ACL 2019]
 - Quantitative knowledge
4. Dice [Chalier et al., AKBC 2020]
 - Multifaceted quantitative contextualization and joint consolidation

WebChild

- Among the first large-scale attempts at text extraction
- Named for getting children's knowledge from the web
- **Focus**: Linking nouns with plausible adjectives
- **Source**: Google web search 5-gram corpus
- **Extraction method**: patterns, ~20 copula verbs (be, look, feel, ...)
- **Contextualization**: Single numeric score
- **Consolidation**: Jointly (label propagation on graph)

Key ideas of WebChild

Volcano is hot.

Chili is hot.

Pop singer is hot.

Text extraction needs semantic refinement

1. Fine-grained relations for commonsense knowledge:

hasAppearance, hasTaste, hasTemperature, hasShape, evokesEmotion,

2. Sense-disambiguation of arguments of knowledge triples (mapped to WordNet):

pop-singer-n¹ hasAppearance hot-a³

chili-n¹ hasTaste hot-a⁹

volcano-n¹ hasTemperature hot-a¹

Approach

For **range and domain population**:

Extract a large list of noisy candidates.

Construct a weighted graph of ambiguous words and their senses.

Mark few seed nodes in the graph.

Use propagation concept: similar nodes (beautiful) (lovely) have similar labels

For **computing assertion**:

Use the range and domain to prune search space of assertions (for a relation)

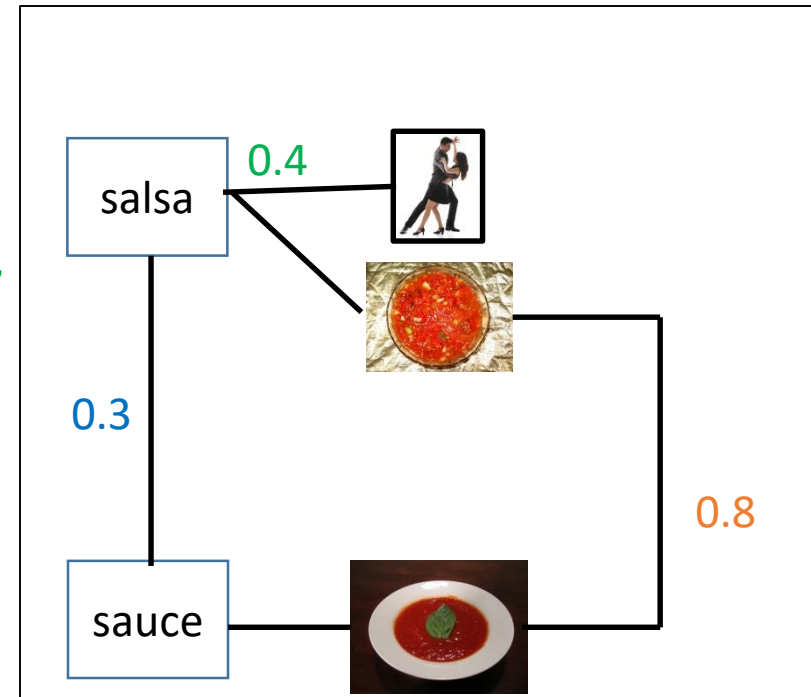
Use propagation concept: similar nodes (car, sweet) (car, lovely) similar labels.

Graph construction per relation (e.g. hasTaste)

- Edge weight:

taxonomic (between senses) ,
co-occurrence statistics (between words),
distributional (between word, senses).

One graph per attr. (here, **hasTaste**)



WebChild: Examples

Domain (hasShape)

face-n¹

leaf-n¹

...

Range (hasShape)

triangular-a¹

tapered-a¹

...

Assertions (hasSshape)

lens-n¹, spherical-a²

palace-n², domed-a¹

...

Sense disambiguation: keyboard-n¹



Top 10 adjectives

ergonomic, foldable, sensitive, black, comfortable, compact, lightweight, comfy, pro, waterproof

Sense disambiguation: keyboard-n²



Top 10 adjectives

universal, magnetic, small, ornamental, decorative, solid, heavy, white, light, cosmetic

Example projects

1. Webchild [Tandon et al., WSDM 2014]
 - Disambiguated noun-adjective pairs
2. **Quasimodo [Romero et al., CIKM 2019]**
 - Salient general triples
3. DoQ [Elazar et al., ACL 2019]
 - Quantitative knowledge
4. Dice [Chalier et al., AKBC 2020]
 - Multifaceted quantitative contextualization and joint consolidation

Quasimodo

= Query Logs and QA Forums for Salient Commonsense Definitions

- Focus on **salient** knowledge
 - Human associations, curiosity
- **Source:** Query logs and QA forum questions
- **Extraction method:** OpenIE
- **Contextualization:** Supervised precision + IDF
- **Consolidation:** Largely per-statement regression



(The Hunchback of Notre Dame)

Starting point: Humans vs. automated IE

Manual constructions:

- Salient but few

[ConceptNet]

elephant is capable of...

 carry a trunk →

 remember water sources →

(6 more)

Automated construction:

- Many but boring

[TupleKB]

Elephant:

- *require, ground*

- *inhabit, region*

- *(95 more)*

How to reconcile the two?

Salient knowledge: Utterance context

Key idea: Questions convey salient knowledge

- Why do cats purr?
- Why do Americans love guns?
- Why are airplanes white?
 - a) So someone knows these!
 - b) That someone cares enough to ask!

Salient knowledge: Premier sources

- QA forums:

- Reddit
- Quora
- Yahoo answers
- Ask.com

- Search engine query logs

- Bing
- Google

Tapping search engine query logs

why do cats

why do cats **purr**

why do cats **like boxes**

why do cats **meow**

why do cats **knead**

why do cats **sleep so much**

why do cats **hate water**

why do cats **like catnip**

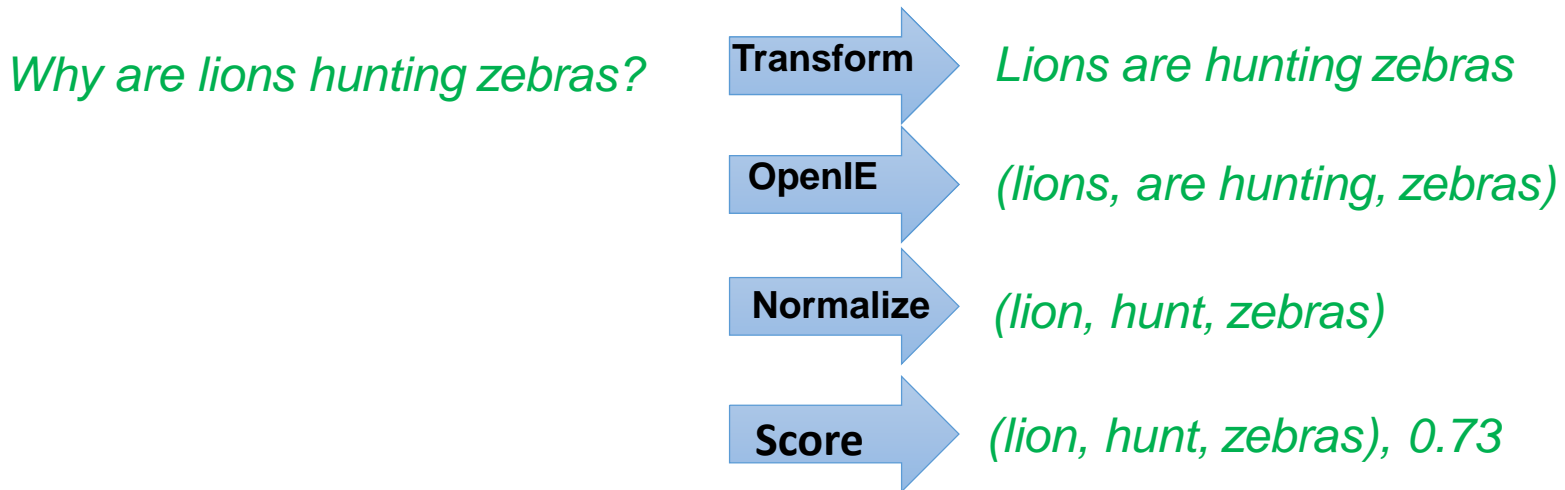
why do cats **lick you**

why do cats **have whiskers**

- Autocomplete gives only 10 suggestions/query
 - Exhaustive suffix probing
 - *Why do cats a*
 - *Why do cats b*
 - *Why do cats ...*
 - *Why do cats aa*
 - *Why do cats ab*
 - ...

Statement extraction

- Questions → statements → tuples using OpenIE



Anecdotal Examples

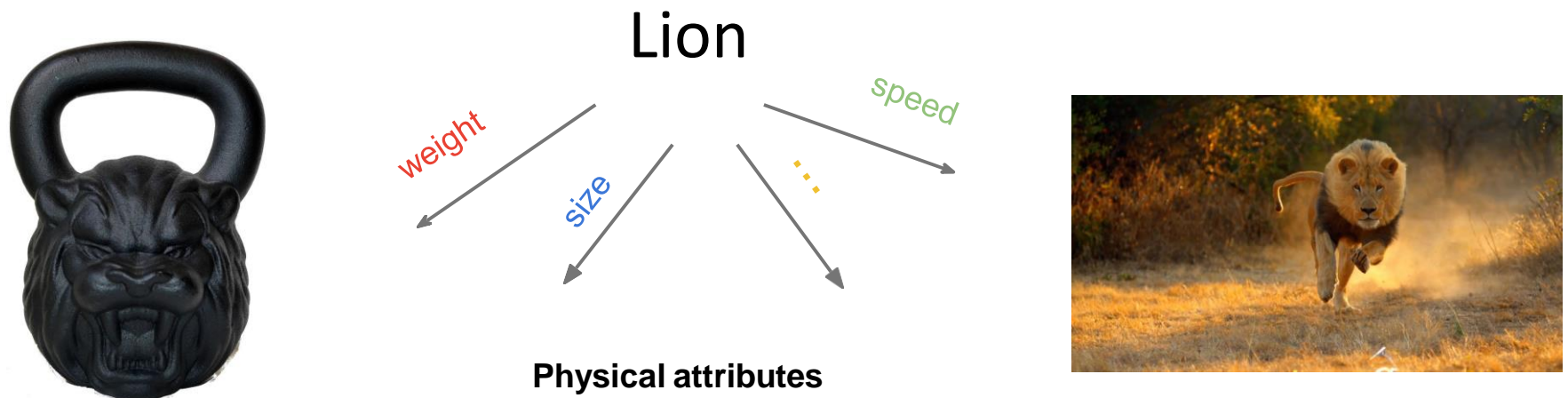
Practical human knowledge	(car, slip on, ice)
Problems linked to a subject	(pen, can, leak)
Emotions linked to events	(divorce, can, hurt)
Human behaviors	(ghost, scare, people)
Visual facts	(road, has_color, black)
Cultural knowledge (USA)	(school, have, locker)
Comparative knowledge	(light, faster than, sound)

Example projects

1. Webchild 1.0 [Tandon et al., WSDM 2014]
 - Disambiguated noun-adjective pairs
2. Quasimodo [Romero et al., CIKM 2019]
 - Salient general triples
3. **DoQ [Elazar et al., ACL 2019]**
 - Quantitative knowledge
4. Dice [Chalier et al., AKBC 2020]
 - Multifaceted quantitative contextualization and joint consolidation

Distribution over quantities (DoQ)

- Understanding numerical properties and the way they relate to words.



- Focus on items which can be measured objectively

Distribution over quantities (DoQ)

- **Source:** Google search engine document index
- **Extraction scheme:** Text window co-occurrence of subject, quantity and dimension keyword
- **Contextualization:** Frequency
- **Consolidation:** none/distribution

Example - Measurement Detection

“These breeds can vary in weight from a

***0.46 kg** teacup poodle ...”*

Detect numerical measurements using rules:

kg/kgs/kilogram -> Mass

Normalize (kg -> g)

Example - Co-Occurring objects

“These ^{Noun} breeds can vary in ^{Noun} weight from a

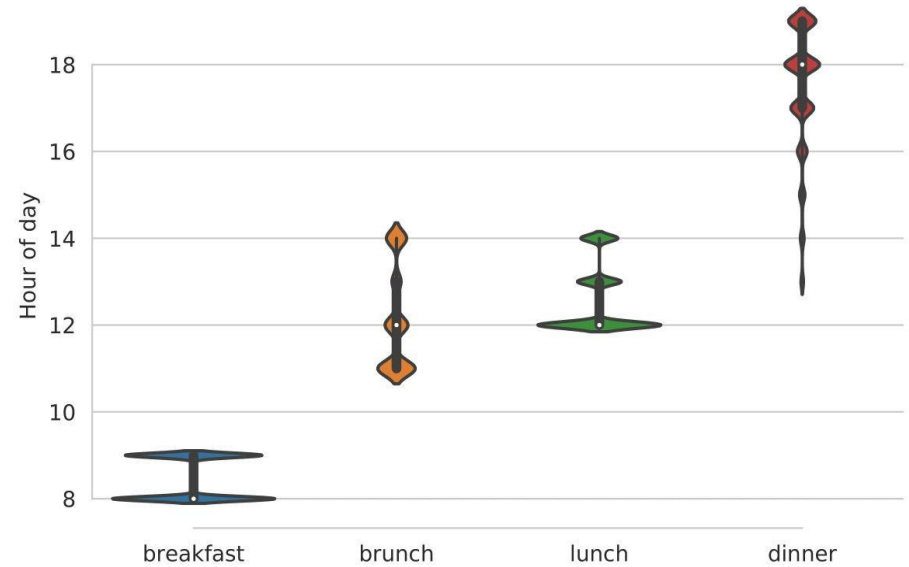
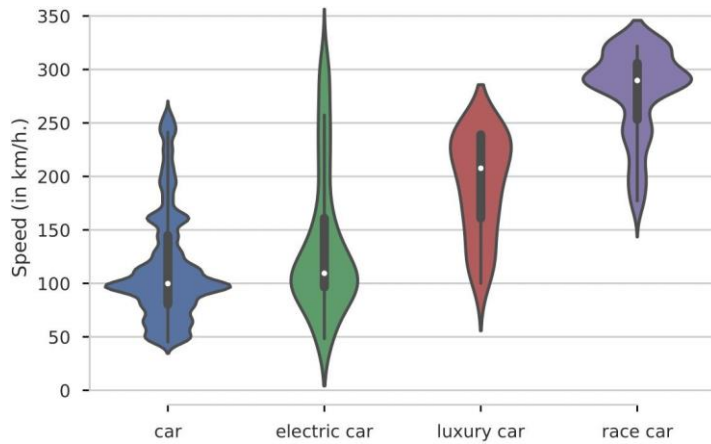
0.46 kg teacup poodle ...”

↙
460 gram

Detect objects of interest (Nouns, Adjectives and Verbs) using a POS tagger.

```
objects_distribution['poodle']['mass'] += [460]
objects_distribution['breeds']['mass'] += [460]
...
```

Example - Aggregating Measurements



Example projects

1. WebChild 1.0 [Tandon et al., WSDM 2014]
 - Disambiguated noun-adjective pairs
2. Quasimodo [Romero et al., CIKM 2019]
 - Salient general triples
3. DoQ [Elazar et al., ACL 2019]
 - Quantitative knowledge
4. **Dice [Chalier et al., AKBC 2020]**
 - Multifaceted quantitative contextualization and joint consolidation

Dice

- A **reasoning framework** for contextualizing existing CSKBs by four numeric facets
 - Plausibility, typicality, remarkability, salience
- **Source**: Any existing CSKB
- **Extraction method**: -
- **Contextualization**: Four numeric facets
- **Consolidation**: Joint taxonomy and similarity-based reasoning

[Chalier et al., AKBC 2020]

A step back – CSK semantics

Lions, attack, humans

A step back – CSK semantics

The semantics we apply to tuples (and which we explain to is true

In WebChild's evaluations

SITUATION:
THERE ARE
14 COMPETING
STANDARDS.

14?! RIDICULOUS!
WE NEED TO DEVELOP
ONE UNIVERSAL STANDARD
THAT COVERS EVERYONE'S
USE CASES.

YEAH!



SOON:
SITUATION:
THERE ARE
15 COMPETING
STANDARDS.

[ConceptNet]

redundancy of terms is captured via inverse document frequency (IDF)
[Information theory 101]

Multi-faceted CSK: Dice

- Each statement (s, p) has **four facets**:

1. Plausibility
2. Typicality
3. Remarkability
4. Salience

- **Lions drink milk** – Plausible, not typical
- **Lions eat meat** – Typical, not salient
- **Lions attack humans** – Salient, not typical

→ **Downstream** tasks left with **all options**



Generic soft constraints for CSK

1. Taxonomical relations give dependencies
 - *Penguins not flying remarkable when most taxonomical siblings do fly*
 - *Macaques eating bananas makes it likely that also stump-tailed macaques eat bananas*
2. Similar statements reinforce each other
 - *Being able to swim correlates with being able to dive*
 - *Lifting logs from the ground correlates with carrying trees*
3. Facets of statements influence each other
 - *Being salient requires being plausible*
 - *Being remarkable and typical implies being salient*

Can combat sparsity!

Can encode coherence expectations!

Dice: Joint reasoning framework

Concept-facets dependencies: $\forall (s, p) \in \mathcal{S} \times \mathcal{P}$

$$\text{Typical}(s, p) \Rightarrow \text{Plausible}(s, p)$$

$$\text{Salient}(s, p) \Rightarrow \text{Plausible}(s, p)$$

$$\text{Typical}(s, p) \wedge \text{Remarkable}(s, p) \Rightarrow \text{Salient}(s, p)$$

Sibling dependencies: $\forall (s_1, p) \in \mathcal{S} \times \mathcal{P}, \forall s_2 \in \text{siblings}(s_1)$

$$\text{Remarkable}(s_1, p) \Rightarrow \neg \text{Remarkable}(s_2, p)$$

$$\text{Typical}(s_1, p) \Rightarrow \neg \text{Remarkable}(s_2, p)$$

$$\neg \text{Plausible}(s_1, p) \wedge \text{Plausible}(s_2, p) \Rightarrow \text{Remarkable}(s_2, p)$$

... parent-child dependencies, similar statement reinforcement

- 17 kinds of soft dependencies in total

Dice: Implementation

Huge constraint system (weighted maxSAT)

How to **bootstrap** constraint system?

- Taxonomy from Hearst-based web extraction [Hertling&Paulheim 2017]
- **Prior scores** from
 - Precision/frequency scores in existing CSKBs,
 - Text entailment models,
 - Statement entropy w.r.t. neighbourhood

How to ground it?

- **Active domain** per subject (+neighbors)
- Still huge **constraint system**
- Approximation via **taxonomy-based slicing**

subject: polar bear

Related concepts

Parents bear, brown bear, mammal, wild animal, predator

Siblings arctic fox, black bear, grizzly bear, panda bear, moose

Facts about 'polar bear'

Click on a property for more details on the statement. Click on a column header to use it as a sorting key.

Show scores as:

Filter by source:

Property	Score	Plausible	Typical	Remarkable	Salient	Source
adapt in summer	0.83	0.19	0.54	0.15	0.15	Quasimodo
adapt to environment	0.83	0.52	0.38	0.93	0.76	Quasimodo
adapt to tundra	0.83	0.10	0.40	0.14	0.10	Quasimodo
be at in arctic	0.67	0.17	0.29	0.93	0.18	ConceptNet
be at risk	0.83	0.62	0.54	0.88	0.80	Quasimodo
be at zoo	0.75	0.10	0.03	0.39	0.37	ConceptNet
be found in arctic	0.91	0.34	0.44	0.51	0.32	Quasimodo
be important to canada	0.92	0.43	0.70	0.27	0.29	Quasimodo
be in danger	0.82	0.91	0.93	0.77	0.97	Quasimodo
be under threat	0.83	0.83	0.80	0.85	0.95	Quasimodo
be used to snow	0.46	0.20	0.51	0.17	0.19	ConceptNet
be white	0.46	0.07	0.68	0.16	0.13	ConceptNet

Outline – Extracting and contextualizing CSK

1. Background
2. Recipe
3. Example projects
4. **Take-away**

Summary

1. Sources

- Domain-specific selection pays off

2. Extraction method

- OpenIE vs. trained extractors

3. Contextualization

- Expressivity-extractability tradeoff
- Quantitative vs. qualitative

4. Consolidation

- Advances in text similarity detection enable joint consolidation

State of the art

- Automatically extracted CSKBs competitive with manually-built projects
 - Usually huge gains in recall, moderate loss in precision

Overview – major projects

	Domain	1. Sources	2. Extraction	3. Contextualization	4. Consolidation	Size (#statements)
WebChild	General noun-adjective pairs	Books	Manual patterns	Single precision	Joint ILP	4.6 M
TupleKB	Science triples	Targeted web search	OpenIE	Single precision	Supervised per-statement	0.3 M
Quasimodo	General triples	User questions	OpenIE	Single precision	Supervised per-statement	4 M (v1.3)
DoQ	Quantity triples	Web crawls	Co-occurrence	Frequency	-	(120 M)
Dice	General triples	Existing structured CSKBs	-	Four quantitative facets	Joint MaxSAT	-
Ascent	General triples	Targeted web search	Facet-based OpenIE	Qualitative facets, subject constraints, frequency	Similarity clustering	8.6 M

Outlook

- Advance of **pre-trained LMs** suggest **hybrid extraction schemes**
 - LMs can contextualize existing uncontextualized CSKBs with plausibility scores
 - Extract salient knowledge directly from LMs
 - Tail knowledge and qualitative contextualizations so far not in reach of pretrained LMs→ See next part
- **Contextualization and ranking of CSK** still open problem
 - Frequency/confidence/plausibility/typicality/salience scores?
 - What kind of qualitative facets?
 - Opportunity for WSDM community

References – Major projects

1. Tandon, Niket, et al. "Webchild: Harvesting and organizing commonsense knowledge from the web." *WSDM*. 2014.
2. Mishra, Bhavana Dalvi, Niket Tandon, and Peter Clark. "Domain-targeted, high precision knowledge extraction." *TACL*. 2017
3. Romero, Julien, et al. "Commonsense properties from query logs and question answering forums." *CIKM*. 2019.
4. Elazar, Yanai, et al. "How large are lions? inducing distributions over quantitative attributes." *ACL*. 2019
5. Chalier, Yohan, et al. "Dice: A Joint Reasoning Framework for Multi-Faceted Commonsense Knowledge" *AKBC*. 2020
6. Nguyen, Tuan-Phong, Simon Razniewski, and Gerhard Weikum. "Advanced Semantics for Commonsense Knowledge Extraction." *WWW* 2021.

Further references

- Omeliyanenko, Janna, et al. "LM4KG: Improving Common Sense Knowledge Graphs with Language Models." *ISWC*, 2020.
- Bhakthavatsalam, Sumithra, Chloe Anastasiades, and Peter Clark. "GenericsKB: A Knowledge Base of Generic Statements." *arXiv preprint arXiv:2005.00660* (2020).
- Bhakthavatsalam, Sumithra, et al. "Do dogs have whiskers? a new knowledge base of haspart relations." *arXiv preprint arXiv:2006.07510* (2020).
- Tandon, Niket, et al. "Knowlywood: Mining activity knowledge from hollywood narratives." *Proceedings of the 24th ACM International on Conference on Information and Knowledge Management*. 2015.
- Chu, Cuong Xuan, Niket Tandon, and Gerhard Weikum. "Distilling task knowledge from how-to communities." *Proceedings of the 26th International Conference on World Wide Web*. 2017.
- Schubert, Lenhart. "Can we derive general world knowledge from texts." *Proc. HLT*. 2002
- Schubert, Lenhart, and Matthew Tong. "Extracting and evaluating general world knowledge from the Brown corpus." *Proceedings of the HLT-NAACL 2003 workshop on Text meaning*. 2003.

Summary

1. **Sources**
 - Domain-specific selection pays off
2. **Extraction method**
 - OpenIE vs. trained extractors
3. **Contextualization**
 - Expressivity-extractability tradeoff
 - Quantitative vs. qualitative
4. **Consolidation**
 - Advances in text similarity detection enable joint consolidation

State of the art

- Automatically extracted CSKBs competitive with manually-built projects
 - Usually huge gains in recall, moderate loss in precision

66

Outlook

- Advance of **pre-trained LMs** suggest **hybrid extraction schemes**
 - LMs can contextualize existing uncontextualized CSKBs with plausibility scores
 - Extract salient knowledge directly from LMs
 - Tail knowledge and qualitative contextualizations so far not in reach of pretrained LMs→ See next part
- **Contextualization of CSK** still with gaps
 - Plausibility vs. typicality vs. salience scores?
 - What kind of qualitative facets?
 - Opportunity for AI community

69

Overview – major projects

	Domain	1. Sources	2. Extraction	3. Contextualization	4. Consolidation	Size
WebChild	General noun-adjective pairs	Books	Manual patterns	Single precision	Joint ILP	4.6 M
TupleKB	Science triples	Targeted web search	OpenIE	Single precision	Supervised per-statement	0.3 M
Quasimodo	General triples	User questions	OpenIE	Single precision	Supervised per-statement	4 M (v1.3)
DoQ	Quantity triples	Web crawls	Co-occurrence	Frequency	-	(120M)
Dice	General triples	Existing structured CSKBs	-	Four quantitative facets	Joint MaxSAT	-
Ascent	General triples	Targeted web search	Facet-based OpenIE	Qualitative facets, subject constraints, frequency	Similarity clustering	8.6 M

67

References – Major projects

1. Tandon, Niket, et al. "Webchild: Harvesting and organizing commonsense knowledge from the web." *WSDM*. 2014.
2. Mishra, Bhavana Dalvi, Niket Tandon, and Peter Clark. "Domain-targeted, high precision knowledge extraction." *TACL*. 2017
3. Romero, Julien, et al. "Commonsense properties from query logs and question answering forums." *CIKM*. 2019.
4. Elazar, Yanai, et al. "How large are lions? inducing distributions over quantitative attributes." *ACL*. 2019
5. Chaliar, Yohan, et al. "Dice: A Joint Reasoning Framework for Multi-Faceted Commonsense Knowledge" *AKBC*. 2020
6. Nguyen, Tuan-Phong, Simon Razniewski, and Gerhard Weikum. "Advanced Semantics for Commonsense Knowledge Extraction." *WWW* 2021.

69

Example projects

1. TupleKB [Mishra et al., TACL 2017]

- Open science triples

2. Ascent [Nguyen et al., WWW 2021]

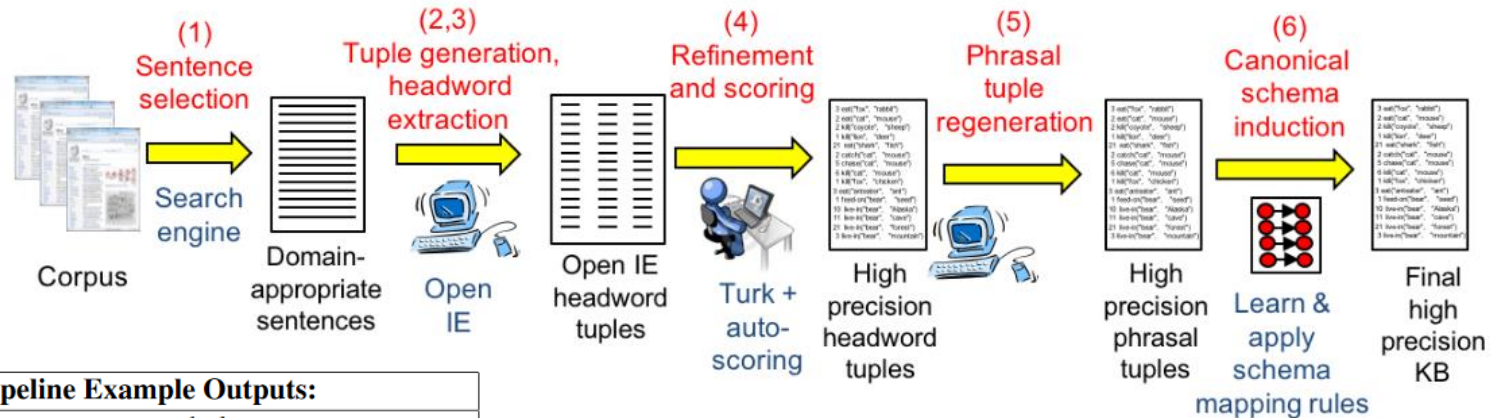
- Qualitative contextualization and state-of-the-art extraction

TupleKB

- Knowledge about **science topics**
- **Source**: Relevant websites via subject-specific keyword queries (template-based)
- **Extraction method**: OpenIE
- **Contextualization**: Single numeric score
- **Consolidation**: Supervised regression per statement

[Mishra et al., TACL 2017]

TupleKB



Pipeline Example Outputs:

Inputs: corpus + vocabulary + types

- 1. Sentence selection:**
 “In addition, green leaves have chlorophyll.”)
- 2. Tuple Generation:**
 (“green leaves” “have” “chlorophyll”)
- 3. Headword Extraction:**
 (“leaf” “have” “chlorophyll”)
- 4. Refinement and Scoring:**
 (“leaf” “have” “chlorophyll”) @0.89 (score)
- 5. Phrasal tuple generation:**
 (“leaf” “have” “chlorophyll”) @0.89 (score)
 (“green leaf” “have” “chlorophyll”) @0.89 (score)
- 6. Relation Canonicalization:**
 (“leaf” “have” “chlorophyll”) @0.89 (score)
 (“green leaf” “have” “chlorophyll”) @0.89 (score)
 (“leaf” “contain” “chlorophyll”) @0.89 (score)
 (“green leaf” “contain” “chlorophyll”) @0.89 (score)

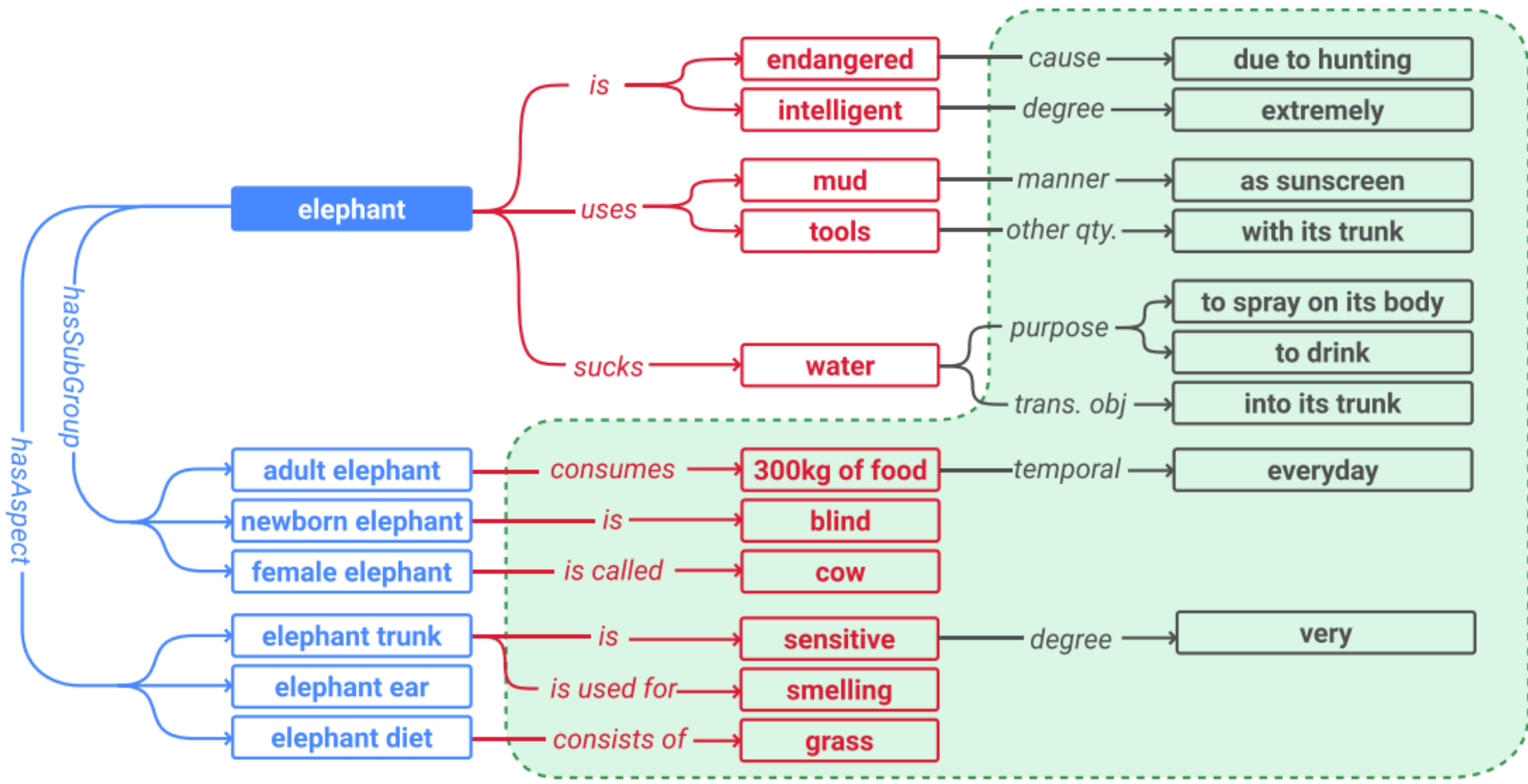
Example projects

1. TupleKB [Mishra et al., TACL 2017]
 - Open science triples
2. **Ascent [Nguyen et al., WWW 2021]**
 - Qualitative contextualization and state-of-the-art extraction

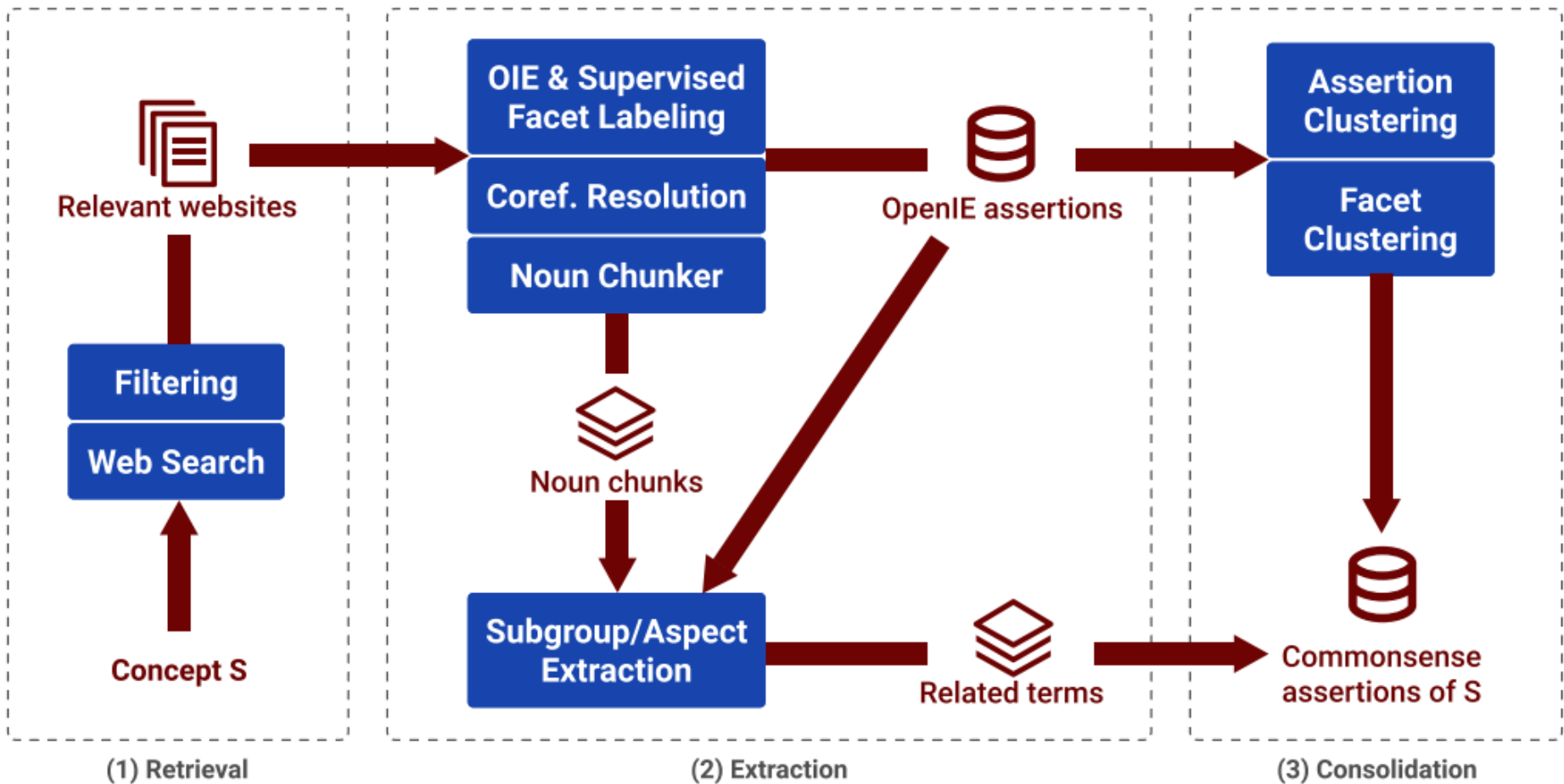
Ascent

- **Source:** Targeted web search
 - Queries created from WordNet hypernyms, e.g., “bank financial institution”
- **Extraction method:** Facet-centric OpenIE
 - Facets give qualitative contextualizations for triples, e.g., location, time, cause, mode
- **Contextualization:** Frequency, qualitative facets, subgroups and aspects
 - *Female elephants, live in, groups, loc: in Africa, 13*
- **Consolidation:** BERT-based clustering

Ascent – qualitative contextualization



Ascent - Architecture



Ascent – BERT-based clustering

Top triple paraphrases

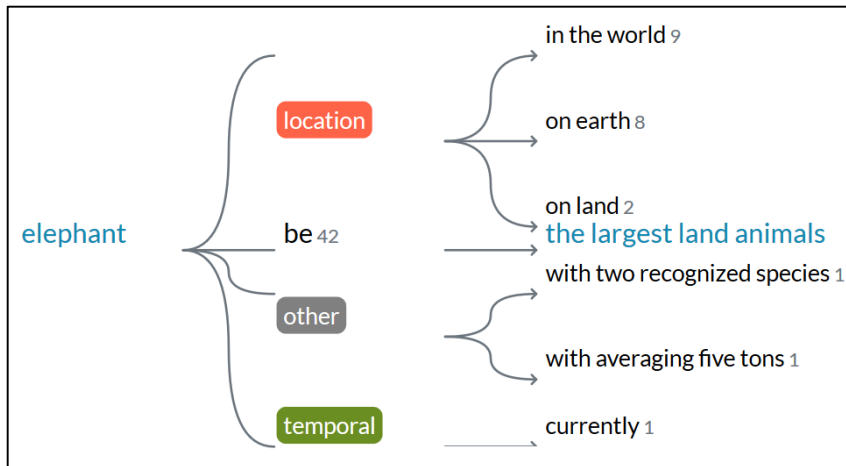
elephant	eat	fruit	13
elephant	consume	fruit	3
elephant	feed on	fruit	2
elephant	feast on	fruit	1
elephant	consume	other fruits	1

▼ 18 source sentences

- Elephants consume grasses, small plants, bushes, **fruit**, twigs, tree bark, and roots.
- Elephants are herbivorous and will eat leaves, twigs, **fruit**, bark and roots.
- Elephants tend to feast on small plants, bushes, **fruit**, twigs, tree bark, and roots and consume up to 330-
- Elephants eat things such as **fruits**, bark, grasses, plants, roots and bushes.
- Elephants are more likely to feed on plants or **fruits** when these are readily available.
- – elephants feed on grass, leaves, tree barks, tender stems and even **fruits**.
- They eat grasses, tree foliage, bark, bamboo, shrubs, roots and **fruit**.
- Elephants eat an extremely varied vegetarian diet, including grass, leaves, twigs, bark, **fruit** and seed pod
- Elephants in Babilie Elephant Sanctuary consume leaves and **fruit** of cherimoya, papaya, banana, guava a
- Elephants are herbivores, consuming ripe bananas, leaves, bamboo, tree bark, and **other fruits**.
- They eat roots, grasses, bark and **fruit** and will even use their tusks to pull off the tree bark or dig in the g

Ascent web interface

Elephant is ...	Elephant has ...	Elephant eats ...
the largest land animals * 42	26 teeth * 8	fruit * 20
herbivore * 35	long trunk 6	grass * 19
intelligent * 32	good memories * 6	plant * 19
endangered * 19	tusk * 6	leaf * 17
good swimmers * 16	four molars * 6	root * 17
more...	more...	more...



<https://ascentkb.herokuapp.com>