# Multi-Cultural Commonsense Knowledge Base Construction
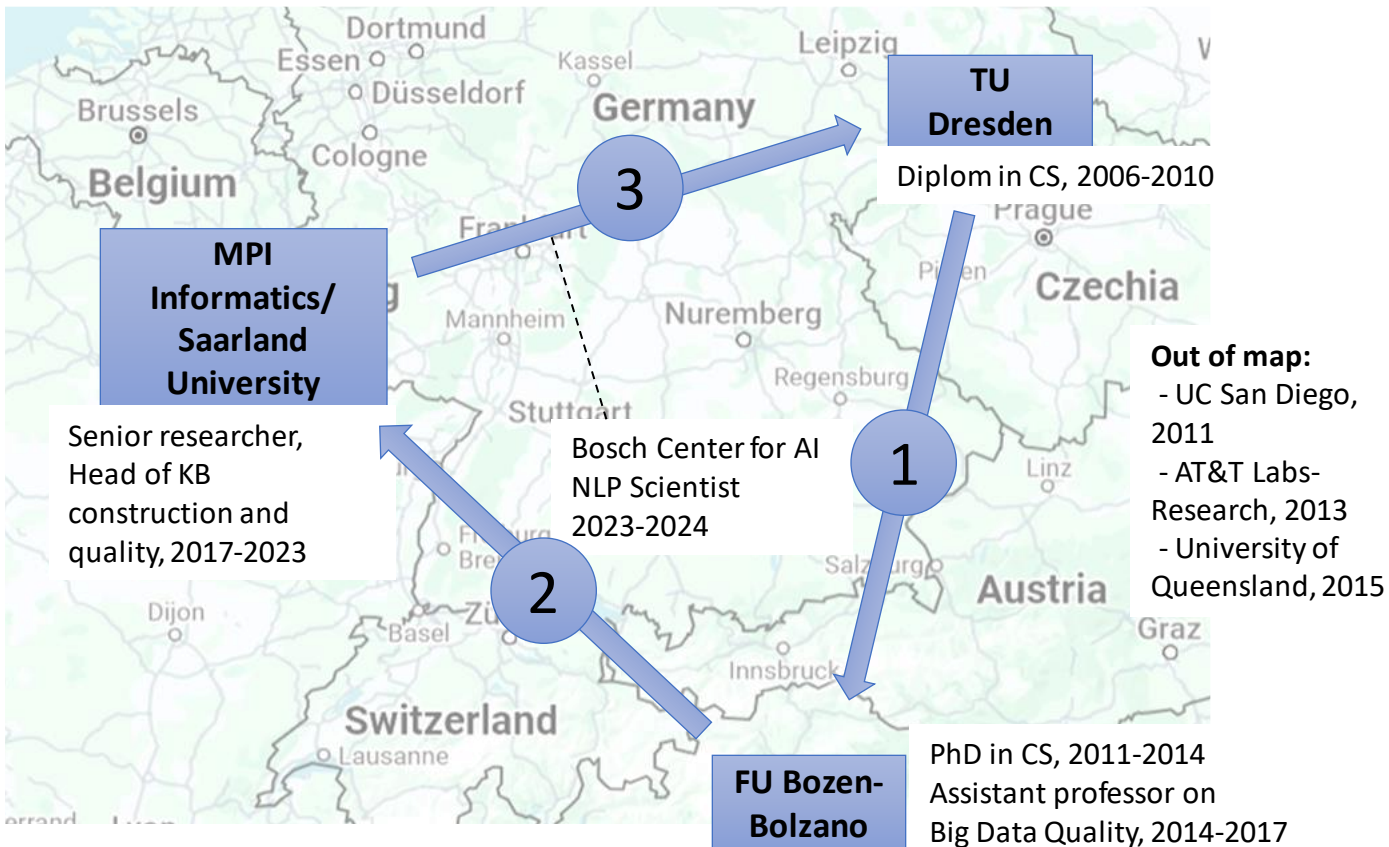
Professorship for Knowledge-aware
Artificial Intelligence

## Simon Razniewski

**TECHNISCHE UNIVERSITÄT DRESDEN**

ScaDS.AI
DRESDEN LEIPZIG

# About myself



TU Dresden
Diplom in CS, 2006-2010

MPI Informatics/ Saarland University
Senior researcher, Head of KB construction and quality, 2017-2023

Bosch Center for AI NLP Scientist 2023-2024

Out of map:
 - UC San Diego, 2011
 - AT&T Labs-Research, 2013
 - University of Queensland, 2015

FU Bozen-Bolzano
PhD in CS, 2011-2014
Assistant professor on Big Data Quality, 2014-2017

# My research agenda

"Develop novel methods for
extracting and consolidating knowledge from, for and with
text, language models (LLMs) and knowledge bases (KBs)"

Current focus:

1. How to know how much KBs/LLMs know?
2. Build high-quality KBs from LLMs
3. Cultural knowledge extraction+consolidation

# My research agenda

"Develop novel methods for
extracting and consolidating knowledge from, for and with
text, language models (LLMs) and knowledge bases (KBs)"

Current focus:

1. **How to know how much KBs/LLMs know?**
2. Build high-quality KBs from LLMs
3. Cultural knowledge extraction+consolidation

# How to know how much KBs/LLMs know?

- Completeness, recall and negation
  - KBs typically operate under OWA
  - When/how can we say that something is not the case (CWA)?
  - Long-standing research line

- Approaches
  - Cardinality assertions in text and KBs    [ACL'17, JWS'20]
    - numberOfSubDivisions(Germany, 16), "consists of 16 states"
  - Conversational maxims classification    [EMNLP'19]
    - "consists of the following states" vs. "the richest states are"
  - Statistical estimators    [WSDM'17]
    - supervised models, species estimation
  - Peer-based inference    [AKBC'20, JWS'21, CIKM'22]
    - you don't have what your peers have

[Razniewski et al., ACM CSUR 2024]    5

# My research agenda

"Develop novel methods for
extracting and consolidating knowledge from, for and with
text, language models (LLMs) and knowledge bases (KBs)"

Current focus:

1. How to know how much KBs/LLMs know?
2. **Build high-quality KBs from LLMs**
3. Cultural knowledge extraction+consolidation

# High-accuracy KBC via LMs

- [Petroni et al., AKBC 2019: Language Models as Knowledge Bases?]  (>2400 citations)

  Prompt: The capital of Saxony is [MASK].

- Limitations
  - Focus on exceedingly popular entities
  - No entity disambiguation
  - Single fact per subject-relation pair
  - Sampling known data → not representative for real use case

# Does this work in practice?

Q1 (Quality): Can we achieve high precision?

- Precision often of utmost importance, e.g., Google KG not deployed because <99% correctness

Q2 (Complementarity):
Can we add value on top of existing KGs?

→Add novel knowledge, not predict existing facts

- Findings (GPT-3):
  - High precision remains tough
  - Best slices: Completion of Wikidata possible at 92% precision on:
    - spokenLanguage, by a factor of ~2.7 (from 2.1Mk to 6.1M)
    - writtenIn, by a factor of ~2.2 (from 14M to 32M)
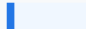    - foundedIn by a factor of ~1.4 (from 43k to 63k)

# Extracting optional and multi-valued relations

- LLMs internally use relative token likelihoods, not truth probabilities
    - →How many objects to retain?
    - →Finding: It is hard

Example: <s> shares a land border with [MASK]

| Prediction | Score |
|---|---|
| Vietnam shares a land border with **Cambodia** . | 12,1 % |
| Vietnam shares a land border with **China** . | 10,7 % |
| Vietnam shares a land border with **India** . | 10,1 % |
| Vietnam shares a land border with **Thailand** . | 9,1 % |

| Prediction | Score |
|---|---|
| Iceland shares a land border with **Norway** . | 18,2 % |
| Iceland shares a land border with **Sweden** . | 7,5 % |
| Iceland shares a land border with **Finland** . | 6,1 % |
| Iceland shares a land border with **Denmark** . | 5,4 % |

# ISWC challenges

- LM-KBC challenge at ISWC 2022
  - Task: Returning ALL correct object values
  - Evaluation on P/R, not just hits@k
  - Sampling mix of popular and long-tail entities

- LM-KBC challenge at ISWC 2023
  - Participants are required to return unique entity identifiers, not just surface form strings

- LM-KBC challenge at ISWC 2024
  - Long lists, numeric attributes, null values frequent

# My research agenda

"Develop novel methods for
extracting and consolidating knowledge from, for and with
text, language models (LLMs) and knowledge bases (KBs)"

Current focus:

1. How to know how much KBs/LLMs know?

2. Build high-quality KBs from LLMs

3. **Cultural knowledge extraction+consolidation**

# Outline

1. **Motivation: Commonsense knowledge**
2. Research challenges
3. Contributions
   A. Representation
   B. Acquisition
   C. Quality assessment
4. Conclusion

# Success of **encyclopedic** knowledge

- Encyclopedic knowledge
  major enabler of knowledge-intensive AI

- Open projects: DBpedia, Wikidata, Yago, etc.
- Proprietary projects at most major tech companies

- Power many applications, e.g., entity disambiguation, question answering, semantic search

- Size: Wikidata: >100M entities, 1.2B statements

# Commonsense knowledge (CSK)

- Statements about classes instead of instances
  - Cities      vs.      Dresden
  - Writers      vs.      Kästner
  - Elephants    vs.      Dumbo



- CSK generalizes, thus more contentious

# Why not simply use **Wikidata**?

<elephant, location, ?>
- Expressible in Wikidata, but no assertions

<lawyer, typical tasks, ?>
- *Give legal advice, represent client, prepare legal documents*
- → Not the typical canonicalized objects
  of encyclopedic KBs

<playing football, prerequisite>
- Subject not even known to encyclopedic KBs

# Why not just LLMs (1)?

- Latent models perform surprisingly well in many tasks

1. But how do they arrive at conclusions?
   →Inherently not scrutable!

2. How can they be modified?
   →No reliable way for adding/removing knowledge

3. What do they actually know?
   → Amount of knowledge not enumerable

# Why not just LLMs (2)?

**Overfitting**     Jabri et al., 2017

Q: Elsa has 5 apples. Anna has 2 more apples than Elsa. _Liz has 4 peaches._ How many apples do they have together?

**Distracted**          Shi et al., 2023

Germans like to meet with friends and family
To the French, family and friends are important
Family and friends are valued highly by the Dutch.

**Overly generic**          [Nguyen et al., 2023]

Who is **Tom Cruise's** mother?

Tom Cruise's mother is **Mary Lee Pfeiffer.**  ✓

Who is **Mary Lee Pfeiffer's** son?

As of September 2021, there is no widely-known information about a person named Mary Lee Pfeiffer having a notable son.  ✗

**Reversal curse**     Berglund et al., 2023

# Why not just LLMs? (3)



→ **Structured knowledge essential to provide a scrutable, reliable, comprehensive data basis for knowledgeable AI**

# CSKB construction: Prior work

- Commonsense reasoning long ambition in AI [McCarthy 1959, Feigenbaum 1984]

- CYC (1980s+): Rich type system with CSK assertions, logical constraints
  - But scaled-down, largely closed-source

- SUMO (2000+): Formal ontology mapped to WordNet

- ConceptNet (1999+)
  - Crowdsourcing for large-scale CSK collection
  - 500k statements of varying quality

- Early LLM prompting (2020-2022)
  - Quality not high enough and overly generic

# Outline

1. Motivation
2. **Research challenges: CSK acquisition**
3. Contributions
   A. Representation
   B. Acquisition
   C. Quality assessment
4. Conclusion -> add ethics

# 1. How should CSK be represented?

- Most common model: <s, p, o> -triples
  - Often with numeric score
- Works from logics and linguistics: Semantic frames
- Works for neural models: Unstructured sentences

→No right or wrong, but where is sweet spot:
 **High expressivity AND sufficient quality data?**

# 2. How can CSK be acquired?

- Previous attempts
  - ConceptNet [1999+]
    - No location for Giraffe
  - WebChild [2014]
    - hyenas are big and small, demonic and fair
  - TupleKB [2017]:
    - <elephant, requires, ground>
    - <elephant, inhabits, region>

Scale

Quality

Salience

- Web is big but full of noise

- **Reconcile scale AND quality?**

Google

"elephants live on the moon"

About 6 results

# 3. How can the quality and coverage of CSK be assessed?

**Intrinsic evaluation**

→ How good is a given statements?

→ Which statements should be acquired?

**Extrinsic evaluation**

→ When is a CSKB useful?

# Outline

1. Motivation
2. Research challenges
3. **Contributions**
   A. **Representation**
   B. Acquisition
   C. Quality assessment
4. Conclusion

# CSK scoring

<Lion, attacks, humans>  - score?
<Lion, drinks, water>     - score?

# CSK scoring

The semantics we apply to tuples (and which we explain to Turkers) is one of plausibility: If the fact is true for some of the arg1's, then score it as true.

[TupleKB]

*In WebChild's evaluations we asked for plausibility*
[WebChild]

| | |
|---|---|
| /r/CapableOf | Something that A can typically do is B. |
| /r/AtLocation | A is a typical location for B, or A is the inherent location of B. Some instances of this would be considered meronyms in WordNet. |
| /r/Causes | A and B are events, and it is typical for A to cause B. |
| /r/LocatedNear | A and B are typically found near each other. Symmetric. |
| /r/Desires | A is a conscious entity that typically wants B. Many assertions of this type use the appropriate language's word for "person" as A. |

[ConceptNet]

*The goal of this paper is to advance the automatic acquisition of salient commonsense properties from online content of the Internet.*
[Quasimodo]

*Informativeness of terms is measured via local frequency and inverse document frequency (TF-IDF)*
[IR theory]

26

# Multi-faceted CSK: Dice

[AKBC'20]



- Each statement has three scores:
  1. Plausibility
  2. Typicality
  3. Salience

- Lion…
  - eats grass – Plausible, not typical
  - drinks water – Typical, not salient
  - attack humans – Salient, not typical

→ Downstream tasks left with all options

# Generic soft constraints for CSK

1. **Taxonomical relations** give dependencies
   - *Lions living in groups salient as most other big cats do not*
   - *<tiger, eats, deer> ⤳ <siberian tiger, eats, deer>*

2. **Similar statements reinforce each other**
   - *<car, causes, accident> ⤳ <car, involved in, crash>*

3. **Facets** of statements **influence each other**
   - *Saliency requires plausibility*
   - *Typicality and frequency imply saliency*

**Deduction rules:**
Can counter sparsity!

**Constraints:**
Can enforce coherence

# Joint reasoning framework

**Scoring-dimension dependencies:** $\forall (s, p) \in \mathcal{S} \times \mathcal{P}$

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

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

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

*¬Typical(Tigers, social) ∧ ¬ Typical(Leopards, social) ∧ … ∧ Typical(Lion, social) ∧ hasParent(Tiger, BigCat) ∧ hasParent(Lion, BigCat), …*
*⊨ Salient(Lion, social)*

… parent-child dependencies, sibling dependencies, similar statement reinforcement

→ 17 types of soft constraints in total

# Joint reasoning: Solution

How to bootstrap constraint system?
- Taxonomy from Hearst-based web extraction [Hertling&Paulheim 2017]
- Prior scores from existing precision/frequency scores, text entailment, entropy

How to ground it?

→ Active domain per subject (+neighbors)

→ Huge constraint system

→ Approximation via taxonomy-based slicing (~100k clauses)

How to solve it?
→Weighted maxSAT
→In general NP-hard
→Constraint shape and linear program approximation make solution tractable (~3 hours @40 cores)

# Triples w/ qualitative facets: Ascent [WWW 2021]

- Quantitative scores often still difficult to interpret
- → Annotate triples with **qualitative** facets
  - Degree, time, location, purpose, instrument, …

  *<elephant, eats, roots;*
  ***degree: sometimes;***
  ***location: forest>***

- Enables higher correctness

  *<elephant, eats, Christmas tree;*
  ***location: zoo>***

- Enables higher informativeness

  *<elephant, uses, their trunk;*
  ***purpose: to suck up water>***

SüddeutscheZeitung

**Weihnachtsbaum-Fütterung der Elefanten im Dresdner Zoo**

11. Januar 2021, 17:42 Uhr

# Cultural CSK: Candle, Mango

- Commonsense knowledge varies widely between cultures and societies
  - Influenced by factors such as geography, religion, occupation

**Candle sample assertions**

| Cultural group | Cultural facet | Cultural commonsense knowledge assertion |
|---|---|---|
| Geo-locations > Countries > Germany | Drinks | German beer festivals in October are a celebration of beer drinking. |
| Geo-locations > Regions > East Asia | Food | Tofu is a major ingredient in many East Asian cuisines. |
| Geo-locations > Regions > South Asia | Traditions | In South Asia, henna is often used in bridal makeup or to celebrate festivals. |
| Occupations > Lawyers | Clothing | Lawyers wear suits to look professional. |
| Occupations > Firefighters | Behaviors | Firefighters run into burning buildings to save lives. |

# Outline

1. Motivation
2. Research challenges
3. **Contributions**
   A. Representation
   B. **Acquisition**
   C. Quality assessment
4. Conclusion

# 2. How can CSK be acquired from online sources?

**Reconcile scale AND quality?**

# Design space

1. Representation: Sentences/triples/frames/…
2. Source: Web/textbooks/Wikipedia/LMs/…
3. Extraction method: Supervised/OpenIE/Prompting…
4. Consolidation: Ranking/clustering/constraint reasoning/…

- Right **design choices** for right **output**:
  - Dice [2020]: Improve quality of existing CSKBs
    - Tuples w/ quantitative facets + existing CSKBs + joint reasoning
  - Ascent [2021]: Extract quality CSK at scale
    - Triples w/ qualitative facets + search engines + OpenIE + clustering
  - Candle [2023]: Empower cultural applications
    - Sentences + web dumps + classification + ranking
  - Also worked with supervised and zero-shot language models (COMET, GPT-3)

→ Experience with a range of sources and techniques

# Web-based CSK acquisition: Ascent

- 10k seed subjects with 500 websites/subject

- Methodology:
  - Template-based query generation
  - Candidate statement extraction via dependency-based patterns (OpenIE)
  - Quality filtering via
    - Targeted disambiguated search engine queries (template/hypernym)

      *Trunk olfactory organ*         *Trunk car part*
    - "Wikipedianess" filter to remove topical outliers

    *Elephants live in …*     ~~*Elephant island is a …*~~     ~~*Elephant Inc. manufactures*~~
  - Semantic grouping by hierarchical clustering

→ Result: 8M statements for 380k subjects

**Elephant**

| | |
|---|---|
| WordNet | elephant.n.01 |
| Wikipedia | Elephant |

**59 salient subgroups of Elephant**

asian elephant 825    african elephant 773    forest elephant 245    bush elephant 181

indian elephant 135    female elephant 133    male elephant 128    more...

**143 salient aspects of Elephant**

trunk 333    tusk 167    ear 166    foot 65    skin 62    mouth 62    teeth 43    more...

**2,828 assertions**

| Elephant is ... | | Elephant has ... | | Elephant is found ... | | Elephant eats ... | |
|---|---|---|---|---|---|---|---|
| the largest land animals * | 44 | 26 teeth * | 8 | in forest * | 9 | grass * | 19 |
| herbivore * | 34 | tusk * | 6 | in desert * | 7 | fruit * | 19 |
| intelligent * | 32 | good memories * | 6 | in africa * | 4 | plant * | 18 |
| endangered * | 22 | long trunk | 6 | in savanna * | 3 | root * | 16 |
| more... | | more... | | more... | | more... | |

1  2  ...  87

**Construction process statistics**

| Bing query | elephant animal facts | Sites crawled successfully | 470 | OpenIE assertions | 50,229 |
|---|---|---|---|---|---|
| Bing results | 500 | Retained sites | 435 | Relevant assertions | 4,085 |
| | | Sentences of retained sites | 28,319 | Clustered assertions | 2,828 |

59 salient subgroups of Elephant



Elephant

| | |
|---|---|
| WordNet | elephant.n |
| Wikipedia | Elephan |

**2,828** assertions

Elephant is ...

the largest land animals
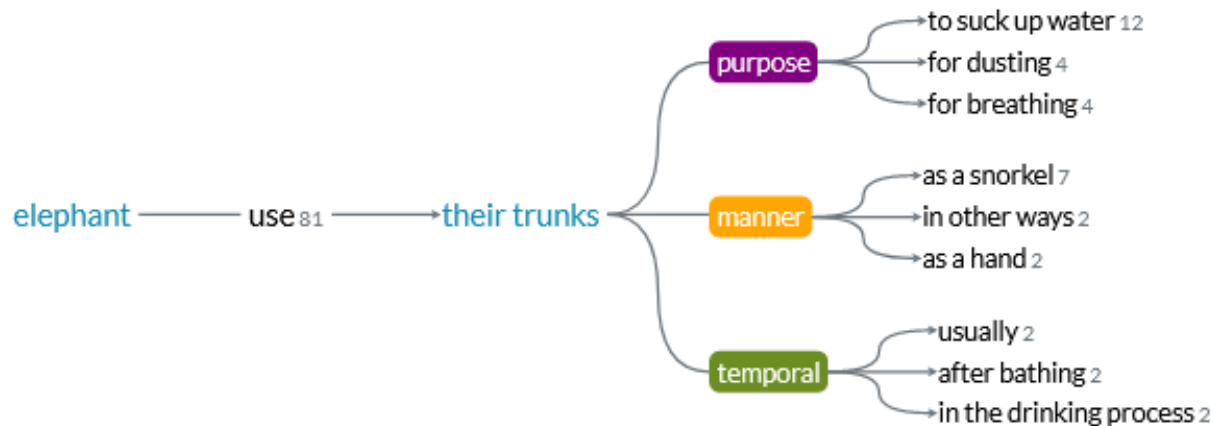
herbivore *

intelligent *

endangered *

more...

Construction process st

| | |
|---|---|
| Bing query | elephant |
| Bing results | 500 |

---

## Assertion summary                                               ✕

elephant ——— use 81 ———→ their trunks

- **purpose**
  - to suck up water 12
  - for dusting 4
  - for breathing 4
- **manner**
  - as a snorkel 7
  - in other ways 2
  - as a hand 2
- **temporal**
  - usually 2
  - after bathing 2
  - in the drinking process 2

### Top triple paraphrases

| | | | |
|---|---|---|---|
| elephant | use | their trunks | **39** |
| elephant | use | its trunk | **23** |
| elephant | use | their trunk | **12** |
| elephant | use | the trunk | **5** |
| elephant | use | trunk | **1** |

19

19

18

16

...

50,229

4,085

38

| | |
|---|---|
| Sentences of retained sites | 28,319 |
| Clustered assertions | 2,828 |

# Cultural CSK: Candle [WWW 2023]

- CSK is conditioned on cultural groups (geography, religion, occupation)

- Methodology
  - Zero-shot language models for classifying relevant sentences
  - Clustering via latent embeddings
  - Structure induction via dictionary-based subject detection and frequency-based concept extraction
  - Ranking via frequency, distinctiveness, specificity, etc.

→ Result: 1.1M sentences
          forming 60k clusters w/ 93k concepts

# Concepts in cultures: Mango

[CIKM 2024]

Can we extract directly from LLMs instead of text?

Result: LLM (chatGPT) yields more AND better quality
- 167k CCSK clusters for 30k concepts and 11k cultures
- Caveat: no text source for GPT-like models openly available

Observations:

1. LLMs can perform induction/interpolation (hallucination)
2. LLM extraction is simpler than text extraction
3. LLM extraction inherently loses source link
   - Especially important in spite of ethical challenges

# Outline

1. Motivation
2. Research challenges
3. **Contributions**
   A. Representation
   B. Acquisition
   C. **Quality assessment**
4. Conclusion

# Intrinsic evaluation

No automated way to assess!
→Need user annotations

Standard metrics (P, R, F1) do not help

→~~Precision~~: Typically several dimensions
  *… is the statement understandable? (meaningfulness)*
  *… would you consider it generally true? (typicality)*
  *… is this common knowledge? (salience)*
  *… does this set the subject apart from others? (distinctiveness)*

→Recall: Evaluated based on similarity-tolerant match to
                                    human associations
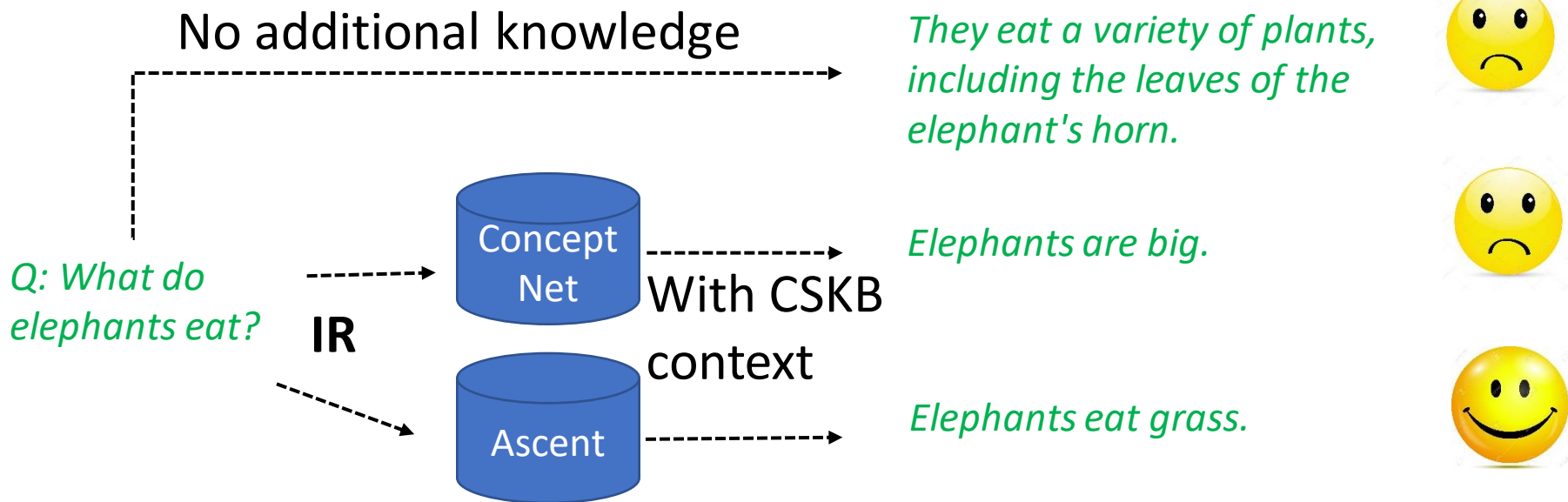
  *Think of lions.*
  *Which 5 statements spontaneously come to mind?*

# Extrinsic evaluation

- Needs a task

- Tasks I have worked with:
  - Multiple-choice question answering
    - Choose answer with most statements connecting it to question
  - Guessing game
    - Given 3 statements, guess the subject
  - Retrieval-augmented generation

# Retrieval augmented generation

**Method**: Serializing KB content for context-enriched LM prompting (BERT, GPT)

No additional knowledge

*They eat a variety of plants, including the leaves of the elephant's horn.*

*Q: What do elephants eat?*

**IR**

Concept Net

With CSKB context

*Elephants are big.*

Ascent

*Elephants eat grass.*

**Result**: Combining Ascent's knowledge with LMs
significantly boosts answer accuracy and informative

**+19% correctness**
**+21% informativeness**

Demo: https://ascent.mpi-inf.mpg.de/qa

44

# Our projects - evaluation

- Intrinsic:
  - Top among automated CSKBs in plausibility/typicality/distinctiveness/cultural relevance
  - Best of all CSKBs in recall

- Extrinsic:
  - Knowledge gives consistent edge in use cases
  - Neural QA models can significantly benefit from symbolic knowledge

# Side node: LLM-CSKBs for LLMs?

- RAG helps LLMs even when the KB is itself generated from an LLM ([CIKM 2024])

- Seemingly cyclic

- Links with insights into chain-of-thought prompting:
  - Giving models "scratch space" before committing to an answer enables them to perform more computations
  - Attention mechanism allows further retrieval

- LLM+CSKBs superior to chain-of-thought in terms of ability to screen knowledge offline

# Outline

1. Motivation
2. Research challenges
3. Contributions
   A. Representation
   B. Acquisition
   C. Quality assessment
4. **Conclusion**

# Do I do logic? (ICCL)

- Symbolic vs. neural AI?
  - Extraction/construction methods use both, with recently more attention to neural approaches
  - Outputs are semi-structured
    - Structure: Concepts, cultural groups, semantic facets
    - Text: Open predicates, sentence assertions
  - Downstream use cases:
    - Commonsense: predominantly RAG (neural),
    - Completeness, encyclopedic KBC: Structured queries (SPARQL)
  - Knowledge acquisition integral for formal reasoning

# Call for connection

- Novel research problems, joint supervision etc.
  - KB-motivated knowledge editing in LLMs
  - A theoretical model for KB evolution
  - LM-KBC at ~Wikidata scale
- Project proposals touching LLMs or KGs, CSK, …
- Scientific event (co-)organization
  - Wikidata workshop?

# Conclusion: Commonsense knowledge

Major challenges in
representation, acquisition, evaluation

Approach:

1. Refined knowledge representation
2. Perform joint reasoning for consolidation
3. Utilize large web excerpts and LLMs with judicious filtering and aggregation

→ **First to combine expressive representations
with large-scale commonsense knowledge acquisition**