

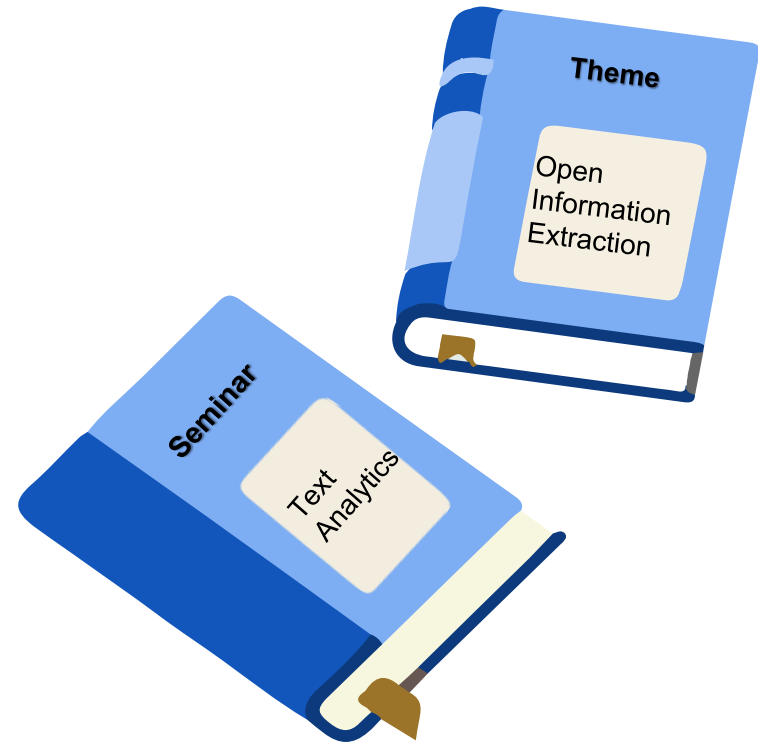
DocOLE : A Document-level Context-Aware Dataset for OpenIE



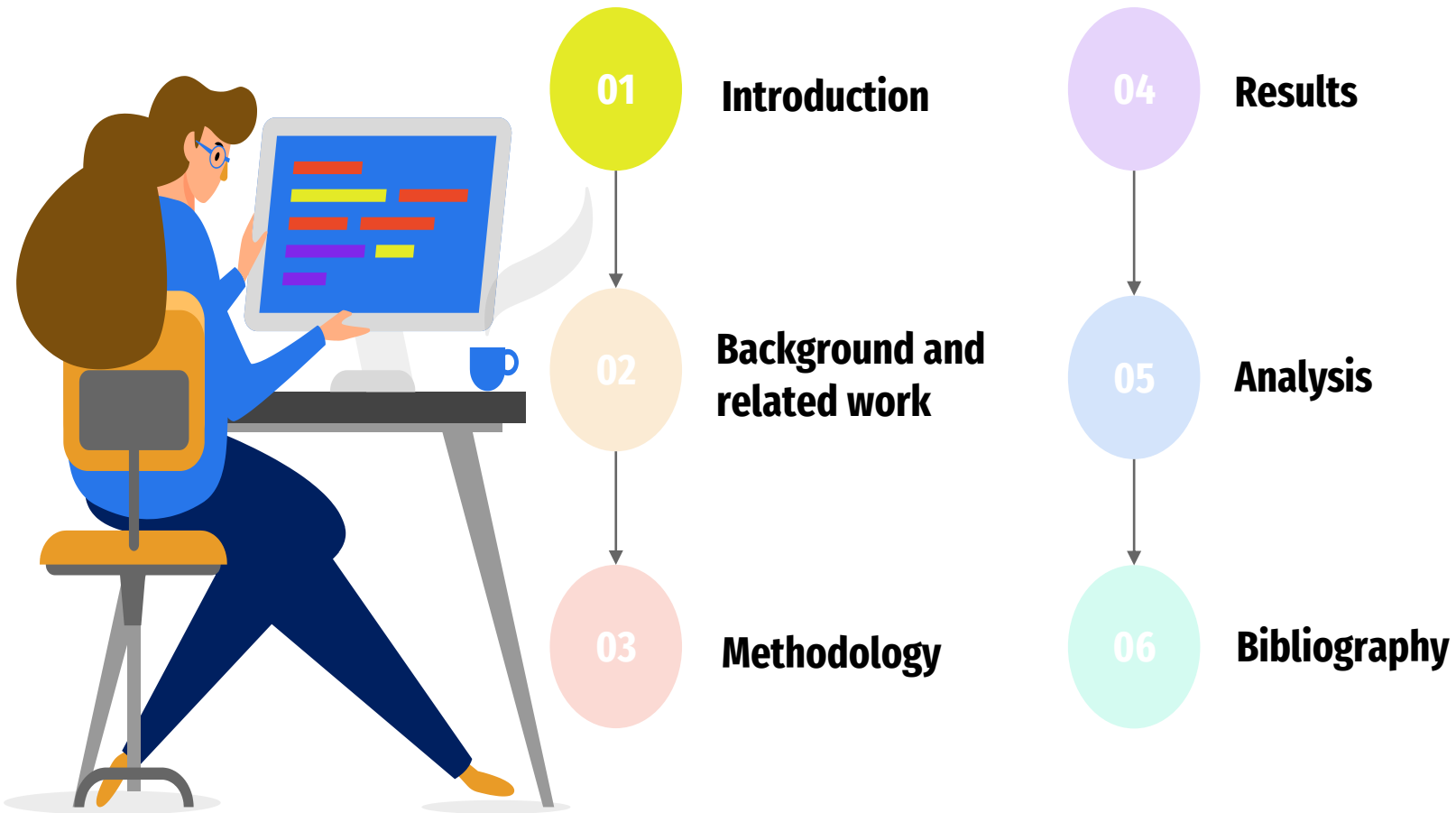
TECHNISCHE
UNIVERSITÄT
DARMSTADT

Kuicai Dong, Yilin Zhao, Aixin Sun, Jung-Jae Kim, Xiaoli Li

Presented by : Amira Chebbi



Overview

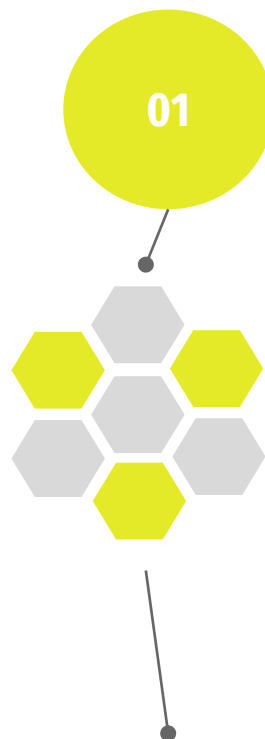


Introduction - Problematic

DocOIE : A Document-level Context-Aware Dataset for OpenIE



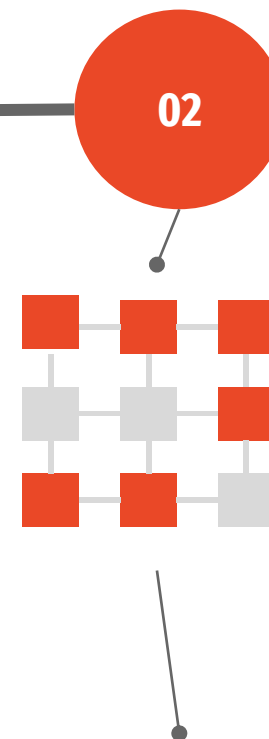
TECHNISCHE
UNIVERSITÄT
DARMSTADT



Part-of-speech Ambiguity

Is it a Verb or a Noun ?

Sentence 1
Data transfers to a single target terminal using the invention might not be significantly faster than conventional download methods. (Pat No. 8495167)
(data; transfers to; a single target terminal) ✗
(data transfers to a single target terminal; use; the invention) ✓
Context S1: Data security is improved as compared with transferring plain text and data transfer requires less time.
Context S2: If a new terminal is registered to the main server during the transfer it will be included in the next data transfer.
Context S3: Examples of data transfers will be described with reference to a preferred embodiment of a network.
Sentence 2
Node-B can be a device a cellular base station having beam-forming antennas that serves various sectors of a cell. (Pat No. 8160027)
(node-B; can be; a device a cellular base station) ✗
(node-B; can be; a device) ✓
(a device, is such as, a cellular base station) ✓
Context S4: A Node-B can be a device such as, a cellular base station that serves an entire cell.



Syntactic Ambiguity

What is the relationship between this two words ?

Introduction - Proposed Solution

DocOIE : A Document-level Context-Aware Dataset for OpenIE



TECHNISCHE
UNIVERSITÄT
DARMSTADT



**DocOIE
proposed
solutions**

01

New Task in OpenIE

Extract relational
tuples with document-
level context

02

DocOIE

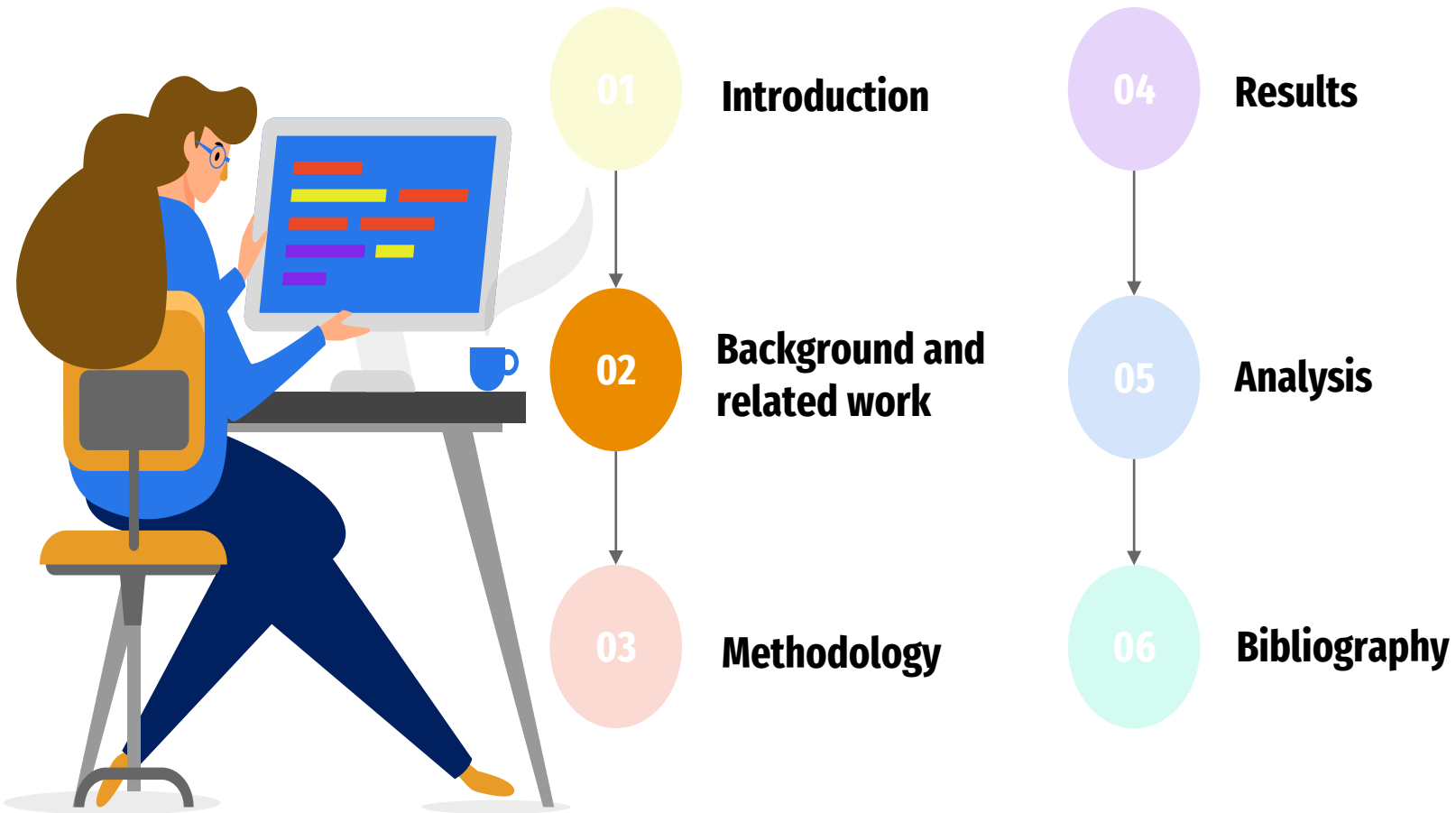
Expert-annotated
dataset

03

DocIE

Neural OpenIE
system

Overview



Related work - OpenIE Datasets

OpenIE Datasets

OIE2016 (Gabriel Stanovsky
and Ido Dagan. 2016)

first large-scale dataset
constructed for OpenIE tasks
and comes with a
standard scoring
framework. the gold tuples are
automatically generated
according to human crafted
rules.

Wire57 (Lechelle et al., 2019)

manually 57
sentences + scoring
improvement

CaRB (Bhardwaj et al. (2019))

50 expert-annotated
sentences + a sophisticated
scoring framework



Related work - OpenIE Datasets - Problems

Dataset's size problem

The number of the annotated sentences remains small

01

No referring to contextual Information

The Annotation is elaborated on a sentence-level and not a document-level

02

OpenIE Datasets Problems



Related works - OpenIE Models

OpenIE Models

TextRunner (Yates et al., 2007)

the first OpenIE system, extract relational tuples based on handcrafted rules or statistical methods

neural OpenIE systems (Cui et al., 2018; Zhan and Zhao, 2020; Kolluru et al., 2020a,b)

They extract tuples in an end-to-end manner, not requiring prior syntactic or semantic analysis.



Related works - OpenIE Models - Problems

Error accumulation for traditional models

Because they rely on prior syntactic or semantic analysis

No consideration of the document-level context

The tuple extraction is based only on the sentence-level

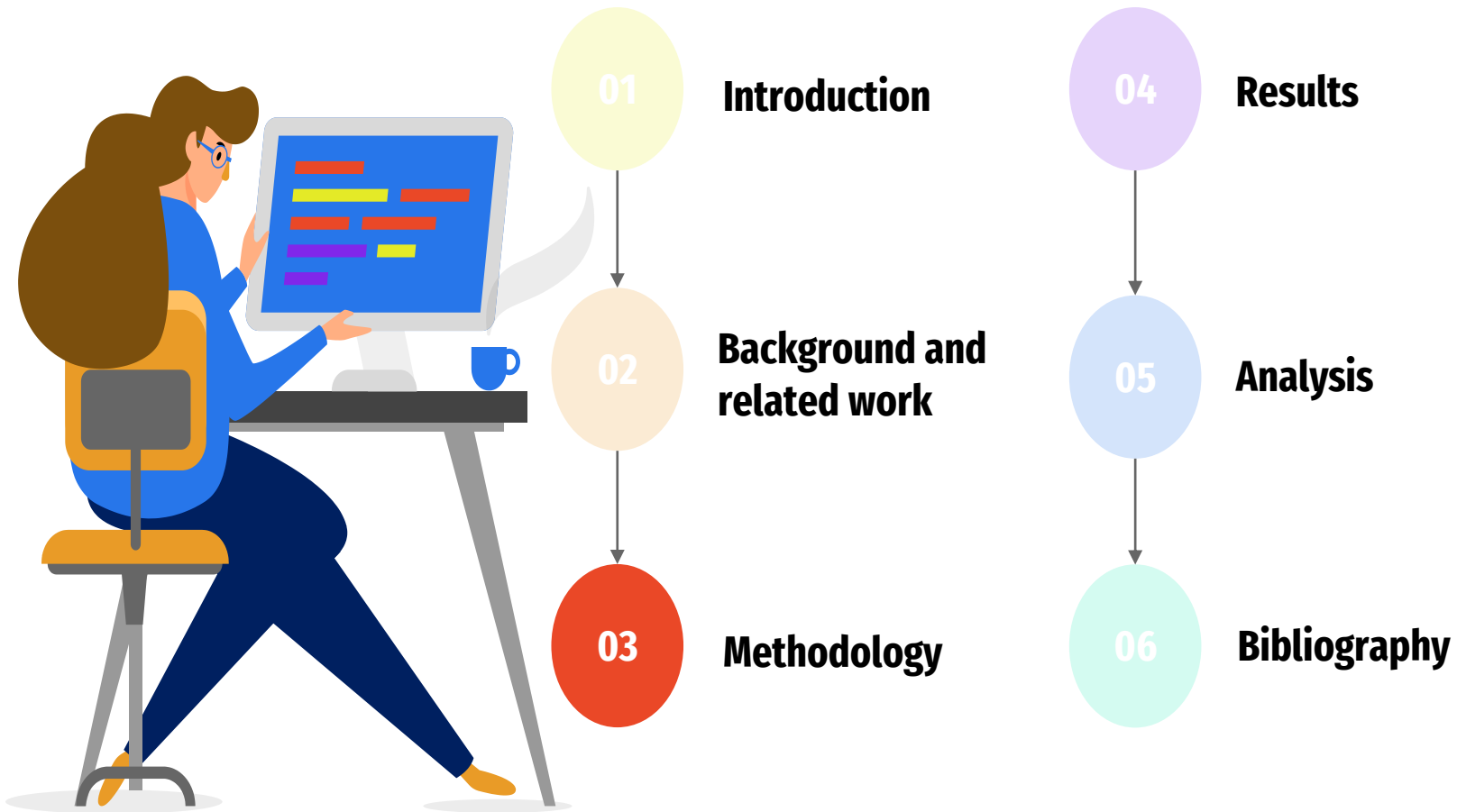
OpenIE Models Problems



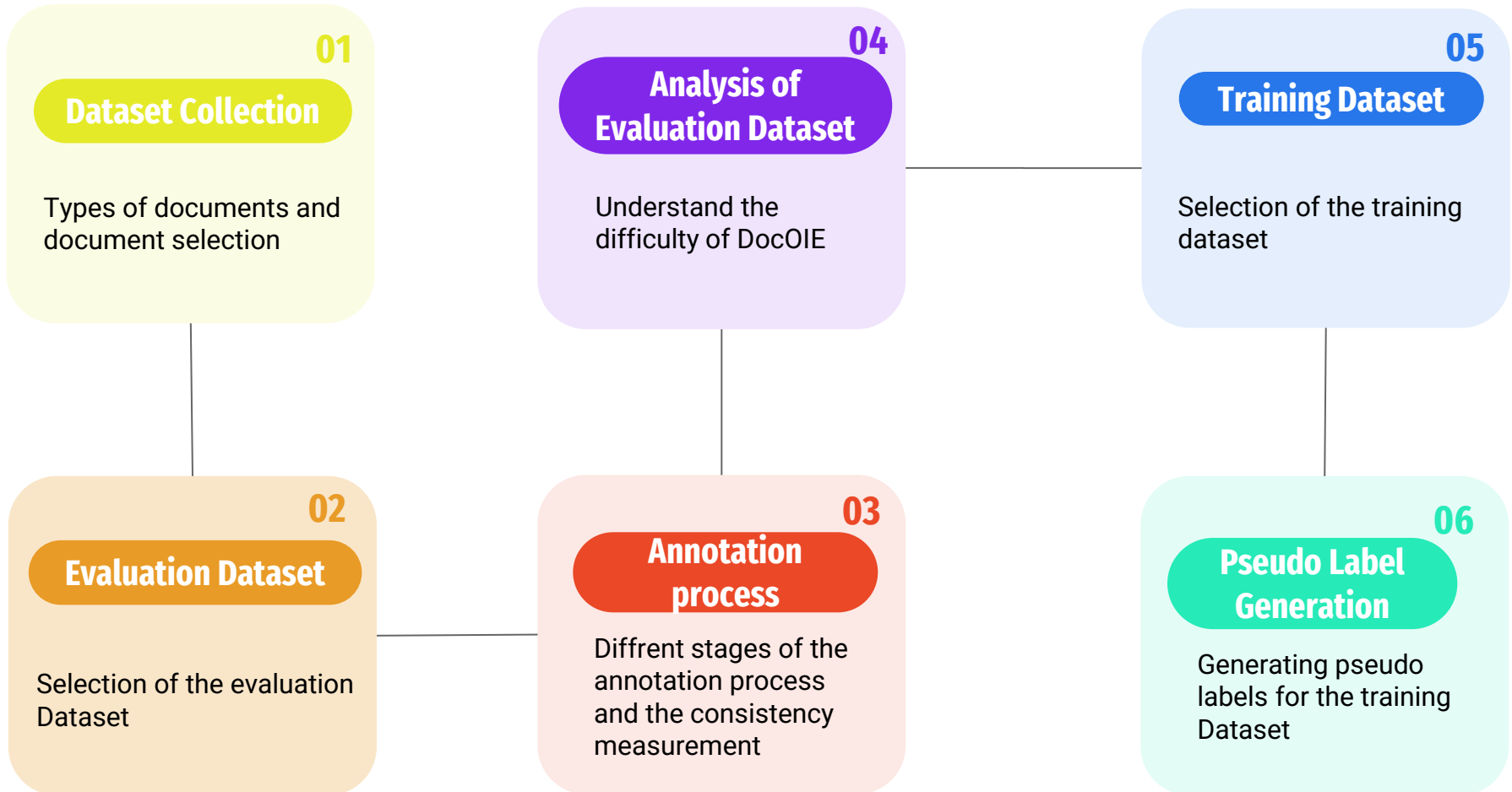
Overview



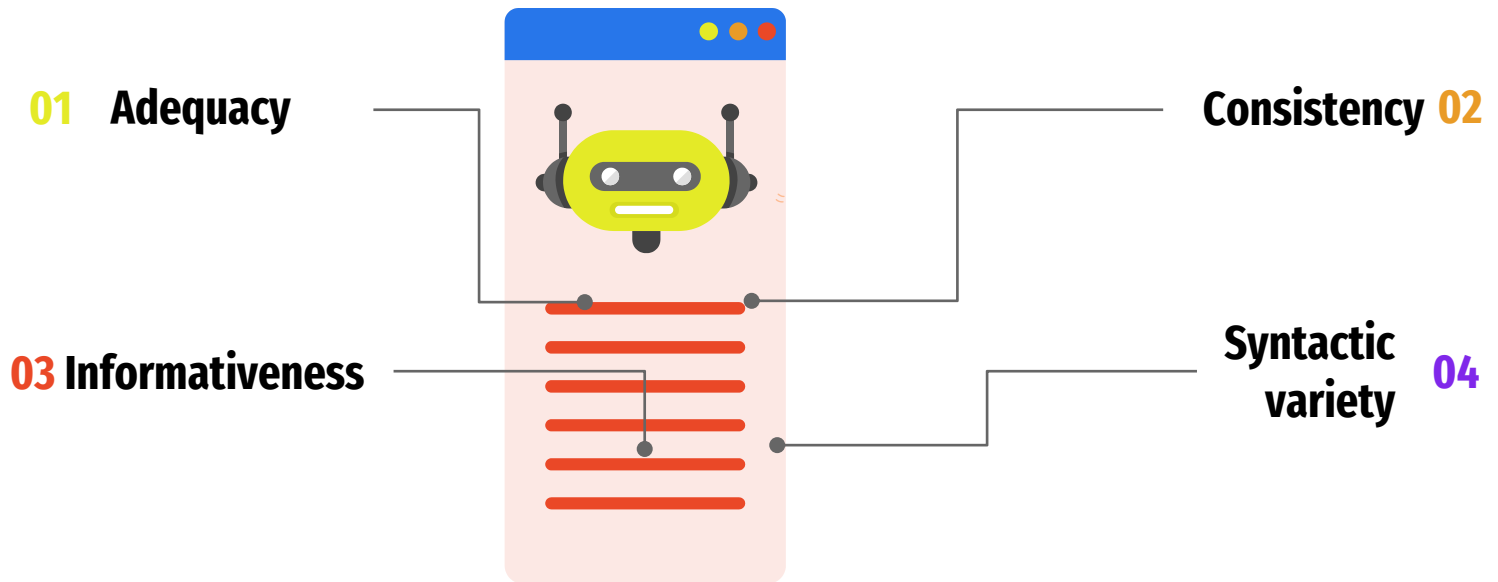
TECHNISCHE
UNIVERSITÄT
DARMSTADT



Methodology – DocOIE Dataset



- Document types : News and Wikipedia articles
- Document selection criterias :



- Keyword selection criterias : Magnitude and Diversity
- Document cleaning

Methodology –DocOIE Dataset

2-Evaluation Dataset

- Randomly select 80 documents (40 each domain)
- From each document , select randomly 10 sentences
- Result : 800 expert-annotated sentences from the DocOIE evaluation Dataset



Annotation process and consistency measurement

Stage 1

1- Annotation

The two annotators practiced annotations independently on 100 sentences

2- Cross validation

Crossvalidation of the results , discussion and then update

Stage 2

1- Annotation

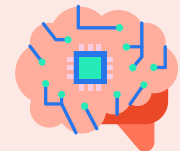
The two experts annotated independently other 100 sentences

2- Consistency Measurement

Evaluation used at tuple-level using the CaRB scorer

Stage 3

Based on the high-level annotation consistency , each expert annotate independently 300 sentences.



Methodology –DocOIE Dataset

4- Analysis of evaluation dataset

- Used to understand the difficulty of DocIE , similar to ([Gashteovski et al., 2019](#)).



Sentence-level Analysis

Evaluate sentence complexity by the number of :

- Conjunction words
- Terminology mention
- Dependent clause

and



Tuple-level Analysis

The analysis of a tuple is based on three points :

- Negative polarity
- Possibility
- Under-specificity

Methodology –DocOIE Dataset

5- Trainig Dataset



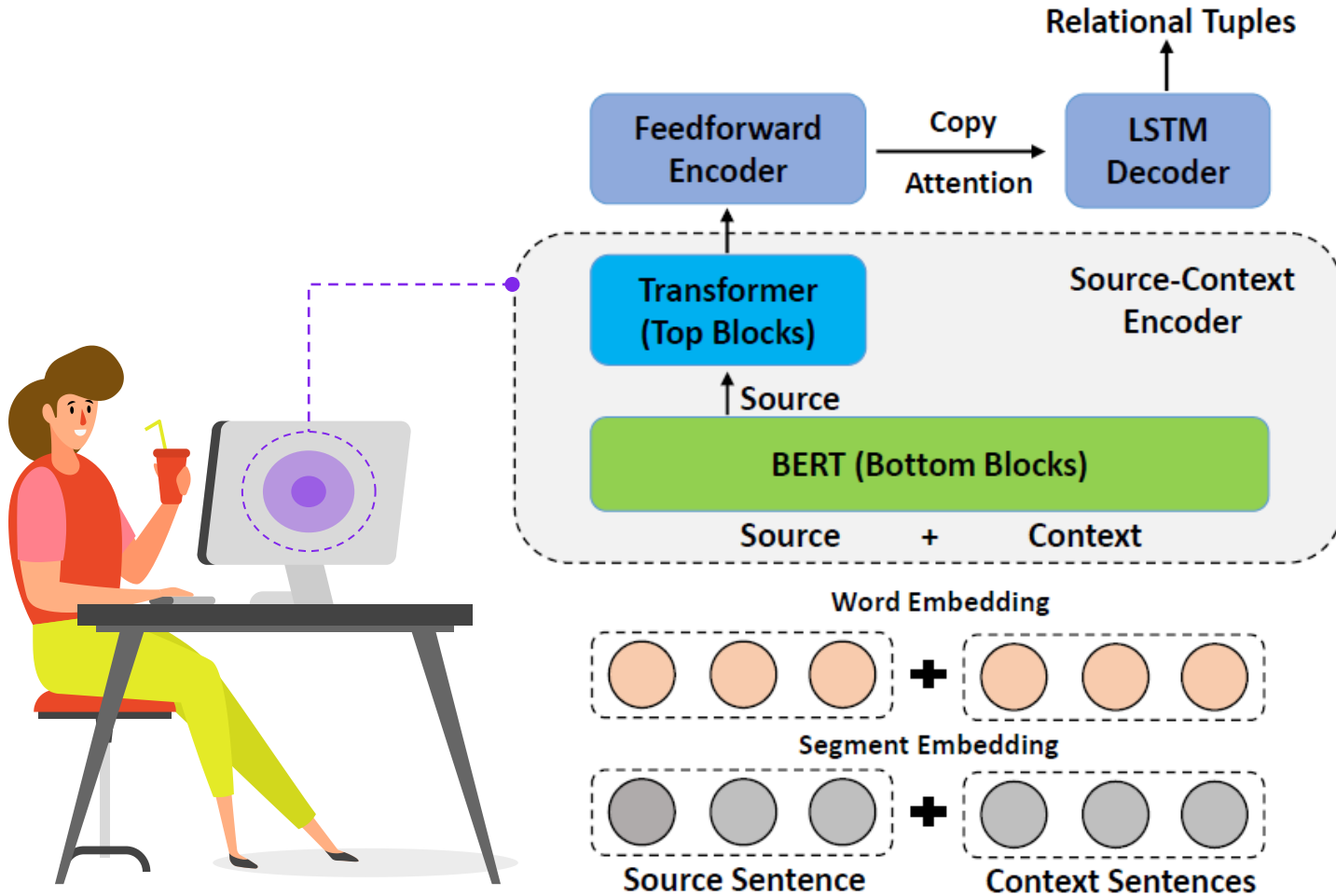
- Select 2400 documents randomly (1200 each Domain)
- The 1200 documents in each domain contain around 120.000 sentences (sufficient for the openIE model)

Methodology –DocOIE Dataset

6- Pseudo Label generation

- Generation of the pseudo labels by bootstrapping with traditional OpenIE models([Kolluru et al.,2020b](#); [Cui et al., 2018](#); [Zhao et al., 2020](#))
- **But First** : Evaluation of the performance of the traditional OpenIE models on the evaluation Dataset using CaRB scorer(to guarantee better quality of pseudo labels)
- The evaluated models were : Reverb ([Fader et al., 2011](#)), Clausie ([Corro and Gemulla, 2013](#)), Stanford OpenIE ([Angeli et al.,2015](#)), OpenIE4 ([Mausam, 2016](#)) , OpenIE5 , Rev+Oie4 and Oie4+ Rev
- Results shows that : both Reverb and OpenIE4 are the best performing individual models and their combinations lead to the best and second best F1 scores in both domains.

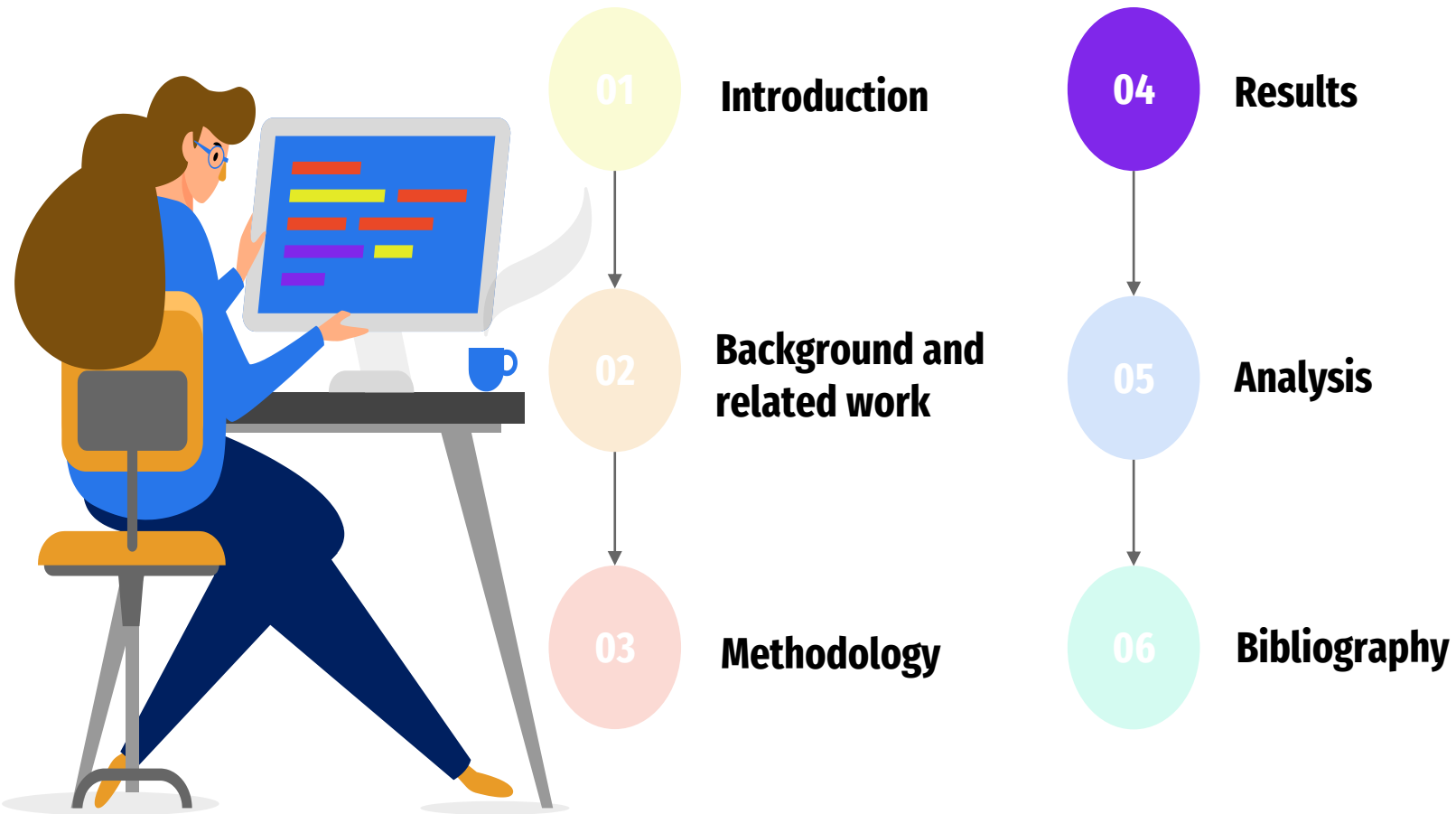
Methodology –DocIE Model



Overview



TECHNISCHE
UNIVERSITÄT
DARMSTADT



Results - DocIE Against sentence-level OpenIE systems

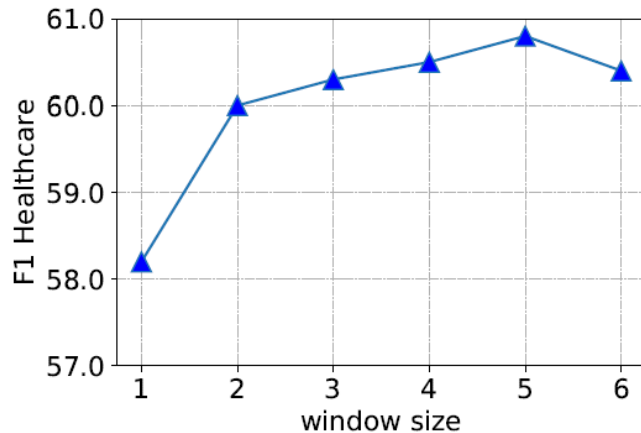


TECHNISCHE
UNIVERSITÄT
DARMSTADT

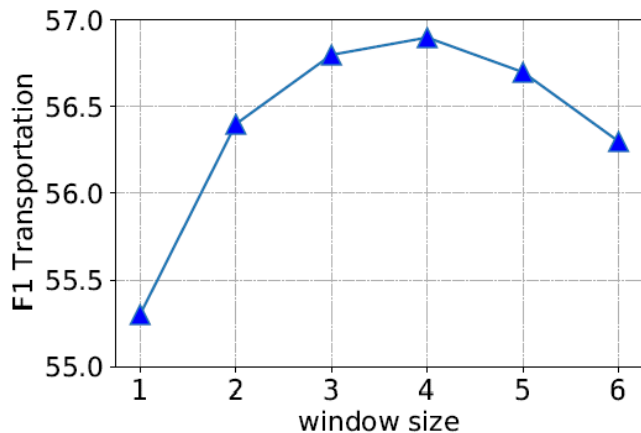
System	Healthcare				Transportation			
	AUC	Prec	Rec	F1	AUC	Prec	Rec	F1
Rev+Oie4	36.8	75.8	47.7	58.6	31.0	74.2	42.4	54.0
Oie4+Rev	35.8	59.6	55.3	57.4	30.1	53.4	52.7	53.0
CopyAttention+BERT	46.8	<u>77.9</u>	48.6	59.8	<u>38.3</u>	55.3	56.9	56.1
IMOJIE	39.7	80.1	46.4	58.7	35.8	<u>63.5</u>	49.2	55.5
DocIE w/o transformer	47.1	76.2	49.9	60.3	38.5	55.8	<u>57.0</u>	56.4
DocIE w transformer	47.4	74.4	<u>51.3</u>	60.8	38.5	56.0	57.5	56.9

- DocIE with transformer achieves the best AUC and F1 in both domains .
- DocIE without transformer is the second best performer and outperforms all the sentence-level systems .

Results - Impact of the context Window size



(a) Healthcare



(b) Transportation

The optimal window size
for the healthcare is 5

The optimal window size
for the transportation is 4

8-10 sentences provide
sufficient context

and

Large window size may
introduce noise



Results - Error Analysis

01 Incompleteness

Fails to cover at least one key phrase in either arguments or relations

02 Incorrect Boundary

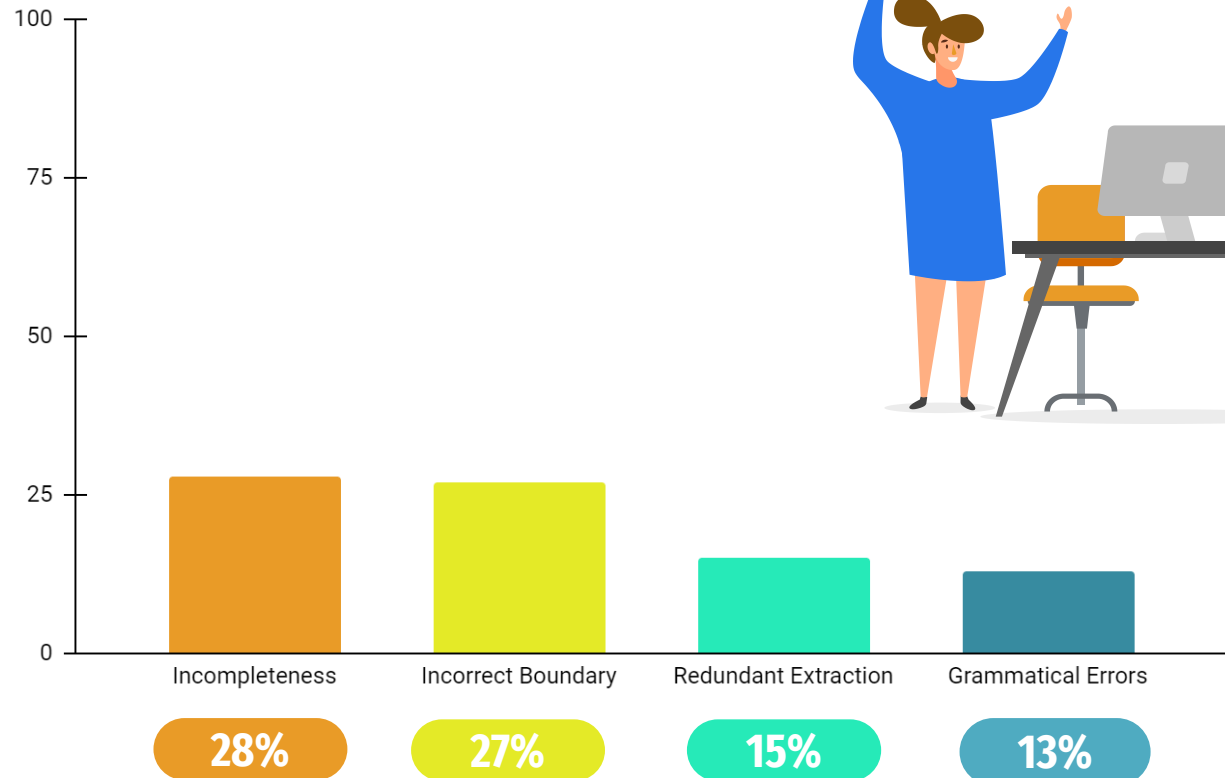
Misinterpretation of the syntactic meaning of the sentence, leading to incorrect Boundary

03 Redundant extraction

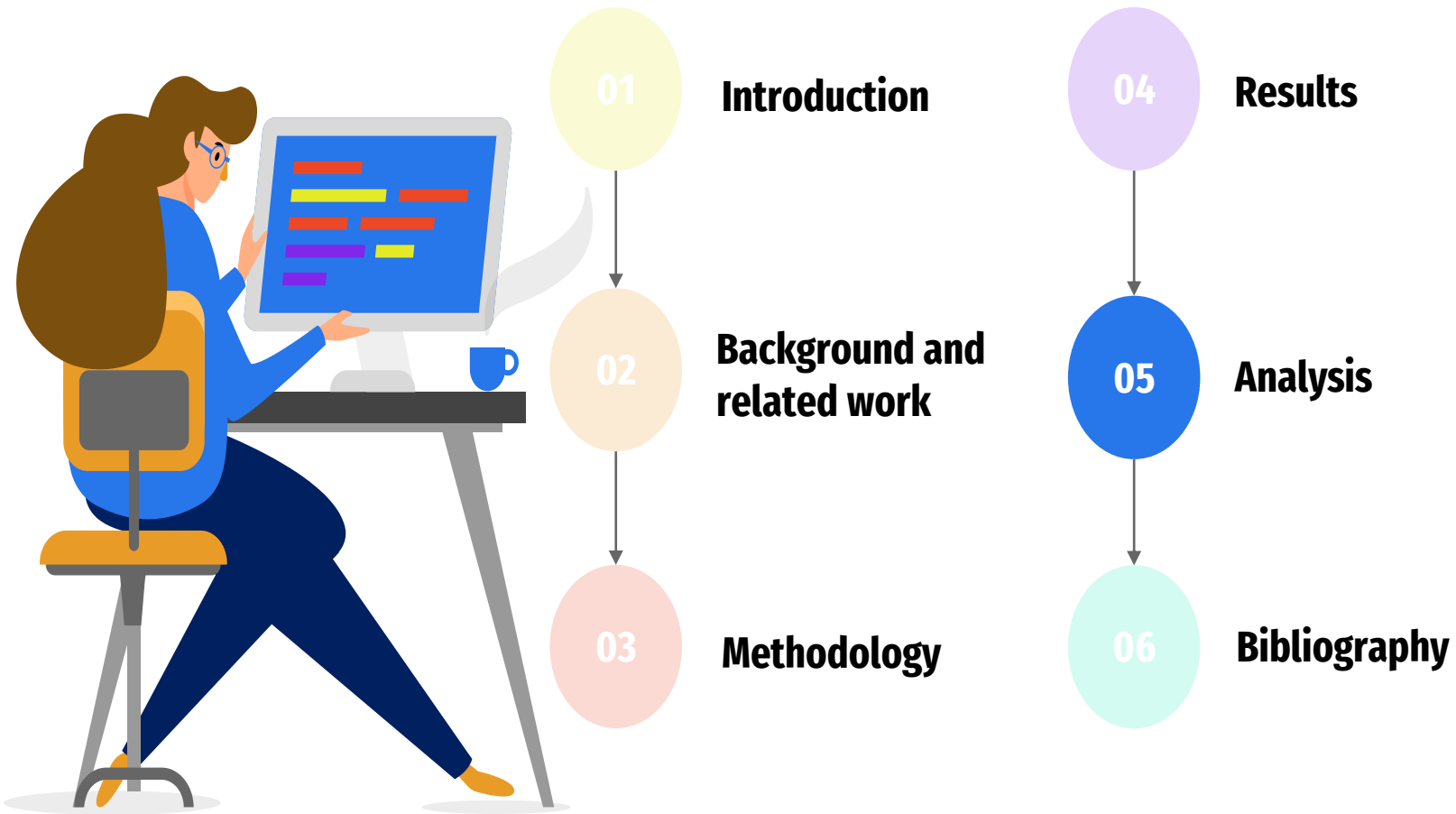
The same relational fact is extracted multiple times

04 Grammatical Errors

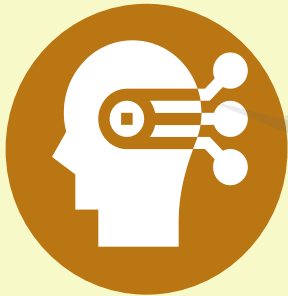
Extractions are not grammatically correct



Overview



Contribution- facet



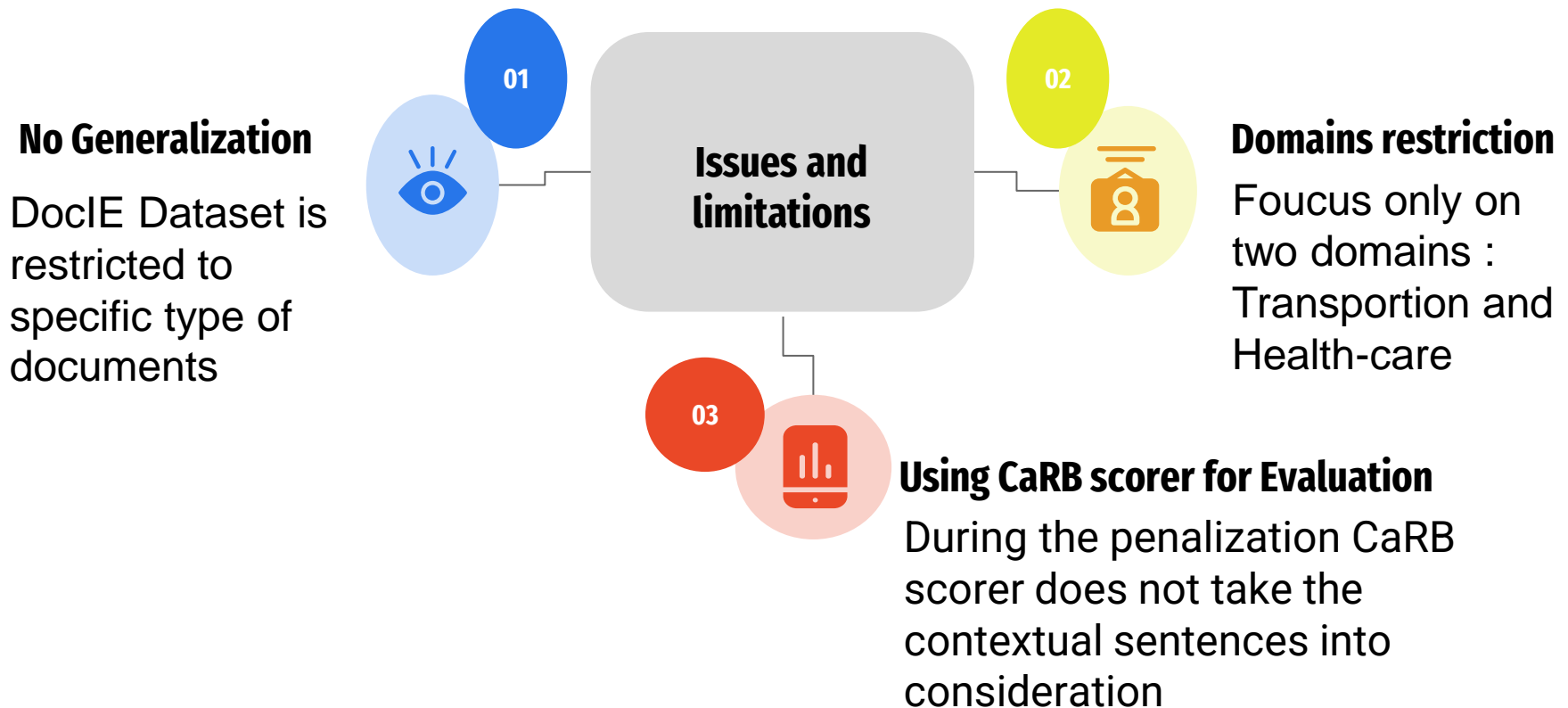
Positive outcomes

Better correct understanding of various topics First document-level context-aware OpenIE dataset and model

Relation to public policy

better analysis and helps in writing the adequate public policy and taking the right measurement in many areas.

Analysis – Possible Issues and limitations



01 Model Improvement

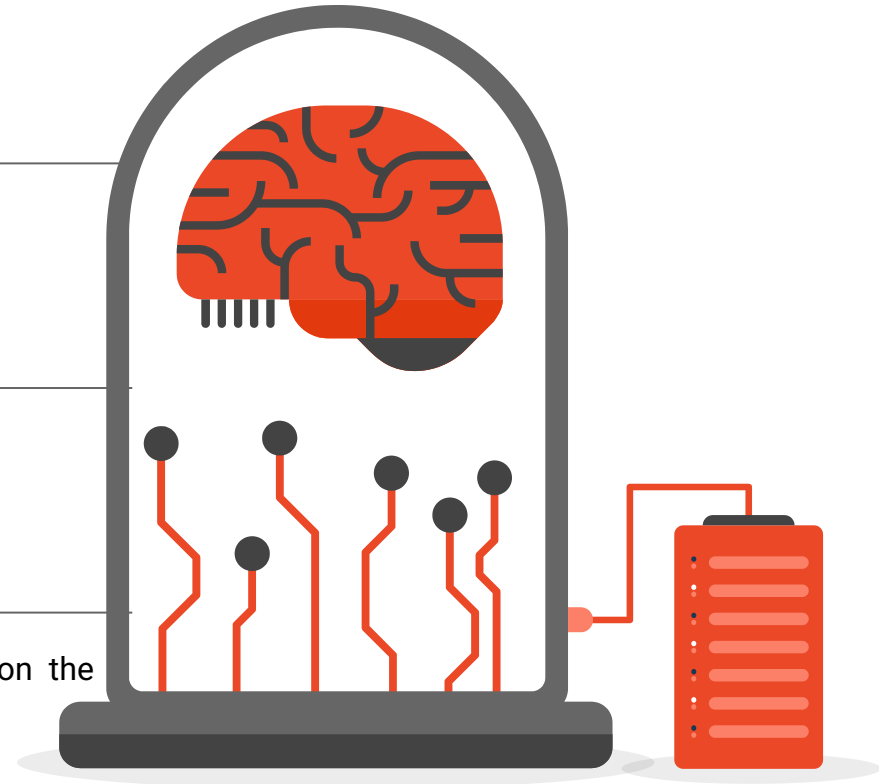
Research on more effective context-aware models

02 Pseudo-Labels

Investigate the possibility of not relying on pseudo labels

03 Adapted Scorer

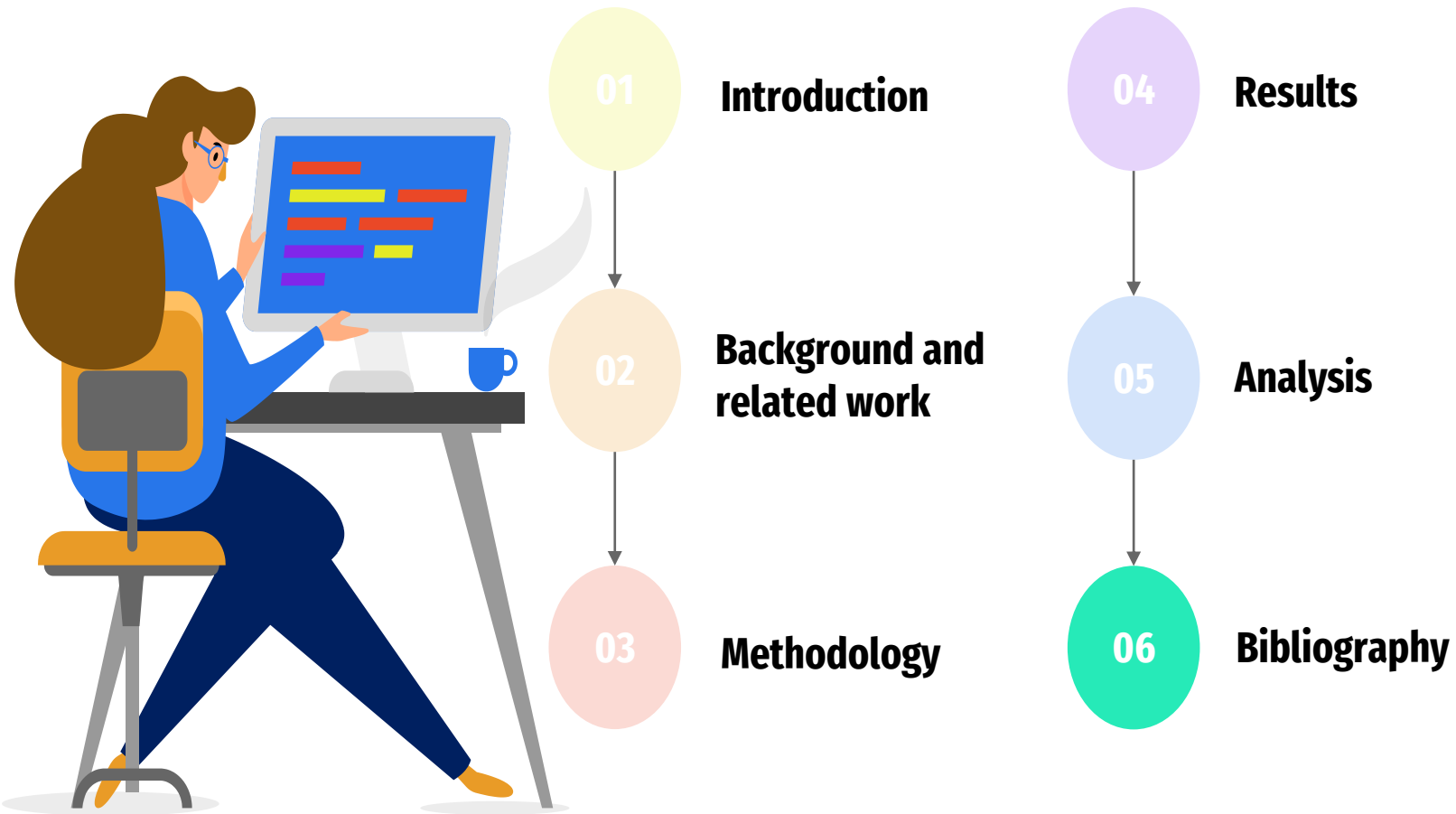
Trying to develop a new scorer that takes into consideration the contextual informations



Overview



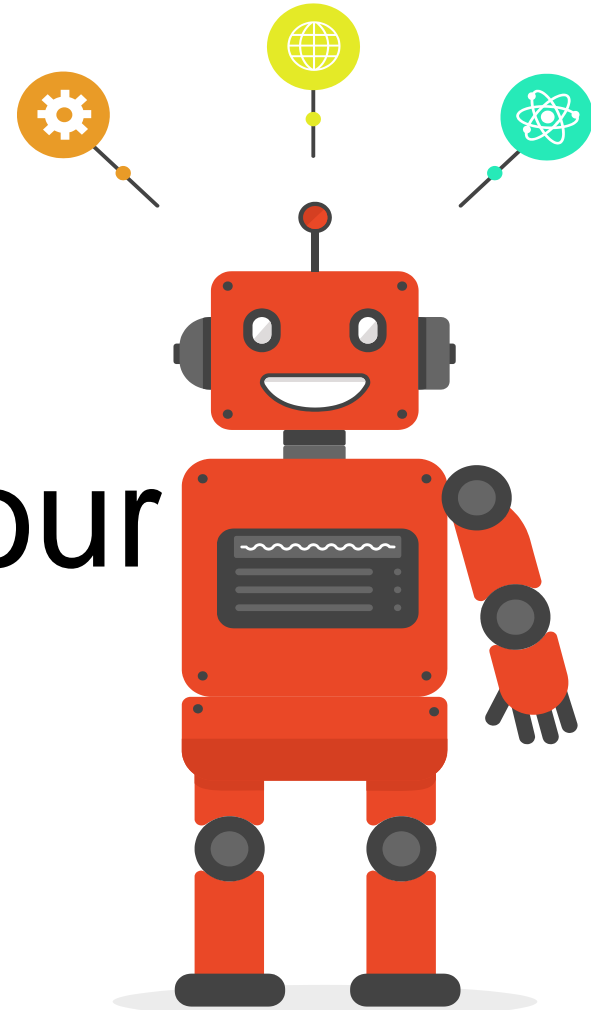
TECHNISCHE
UNIVERSITÄT
DARMSTADT



- DocOIE: A Document-level Context-Aware Dataset for OpenIE
- Creating a Large Benchmark for Open Information Extraction
- WiRe57 : A Fine-Grained Benchmark for Open Information Extraction
- CaRB: A Crowdsourced Benchmark for Open IE
- TextRunner: Open Information Extraction on the Web
- Span Model for Open Information Extraction on Accurate Corpus



Thank you for your Attention

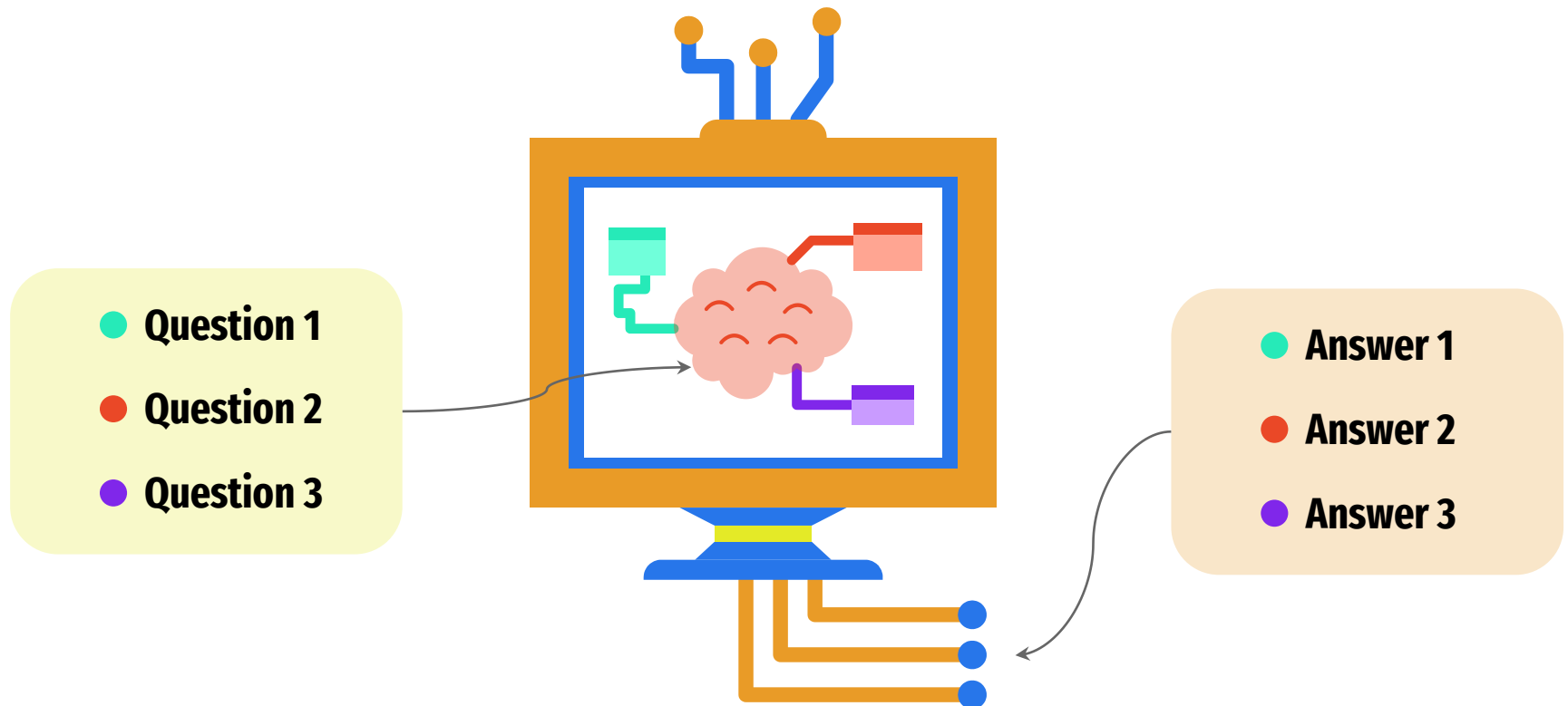


Q/A Part



TECHNISCHE
UNIVERSITÄT
DARMSTADT

Questions



Answers