

NLP Tasks

Natural Language Processing

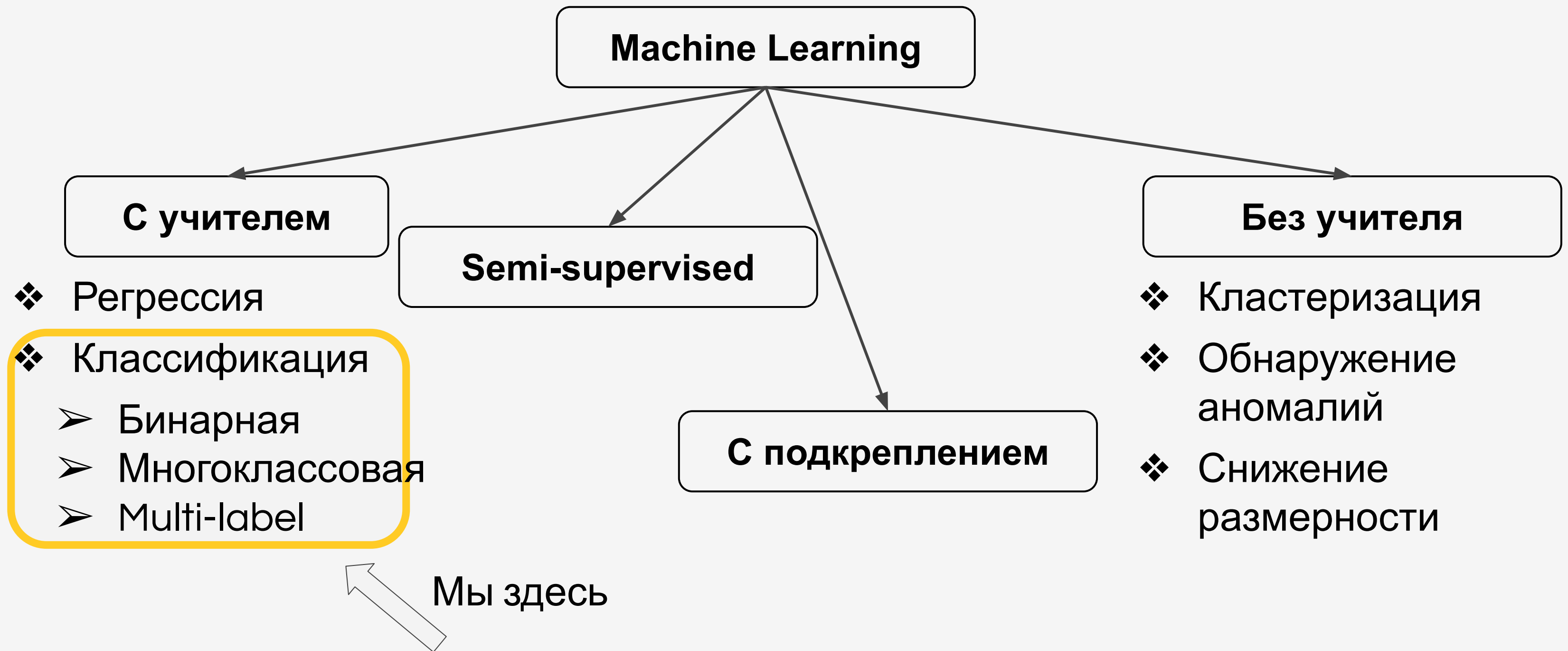
Ника Зыкова, 2025/12/11

План на сегодня

- ❖ Классификация задач:
 - По типу входных и выходных данных;
 - По уровню языка;
 - По используемым методам;
- ❖ Классические NLP задачи:
 - Постановка;
 - Входные/выходные данные;
 - Метрики;
 - Методы.



ML-задачи



NLP-задачи

Token-level

каша

была

вкусная

но

оформление

ужасное

вернусь

еще

NOUN

VERB

ADJ

CONJ

NOUN

ADJ

VERB

ADV

Span-level

Clause 1

Clause 2

Clause 3

Sequence-level

Negative

NLP-задачи

Seq2Seq

каша

была

вкусная

но

оформление

ужасное

вернусь

еще

the porridge

was

delicious

but

the decoration

was terrible

i will be back

again

Auto-regressive

если

исправят

подачу

в

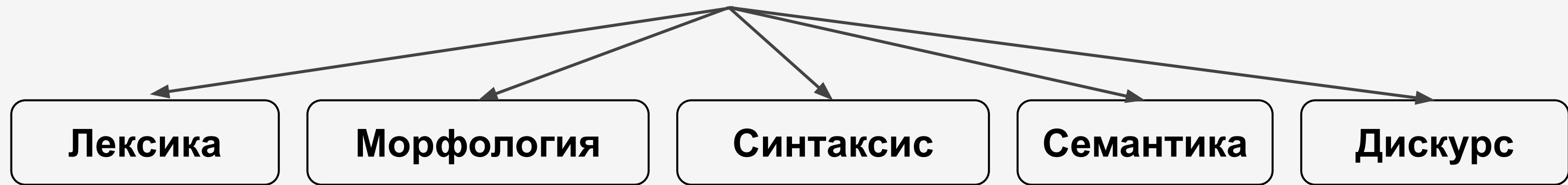
следующий

раз

Structured prediction



По уровню языка

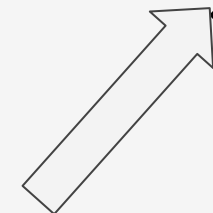


- ❖ Лемматизация
- ❖ Части речи
- ❖ Морфологические характеристики

- ❖ Анализ тональности
- и
- ❖ NER
- ❖ WSI / WSD

- ❖ Кореференция
- ❖ Анафора
- ❖ Диалоги

Почти все
задачи тут



Используемые подходы

- ❖ Rule-based: лексиконы, правилковые парсеры,
- ❖ Статистические методы: n-gram, марковские цепи, CRF;
- ❖ Классический ML: статистические вектора + логрег / случайный лес etc, kernel methods;
- ❖ Нейросетевые: RNN, LSTM, CNN, word2vec и похожие эмбединги;
- ❖ Transformer-based: BERT & Co, своя модель на каждую задачу;
- ❖ LLMs: одна большая языковая модель с широким набором навыков на множество задач.

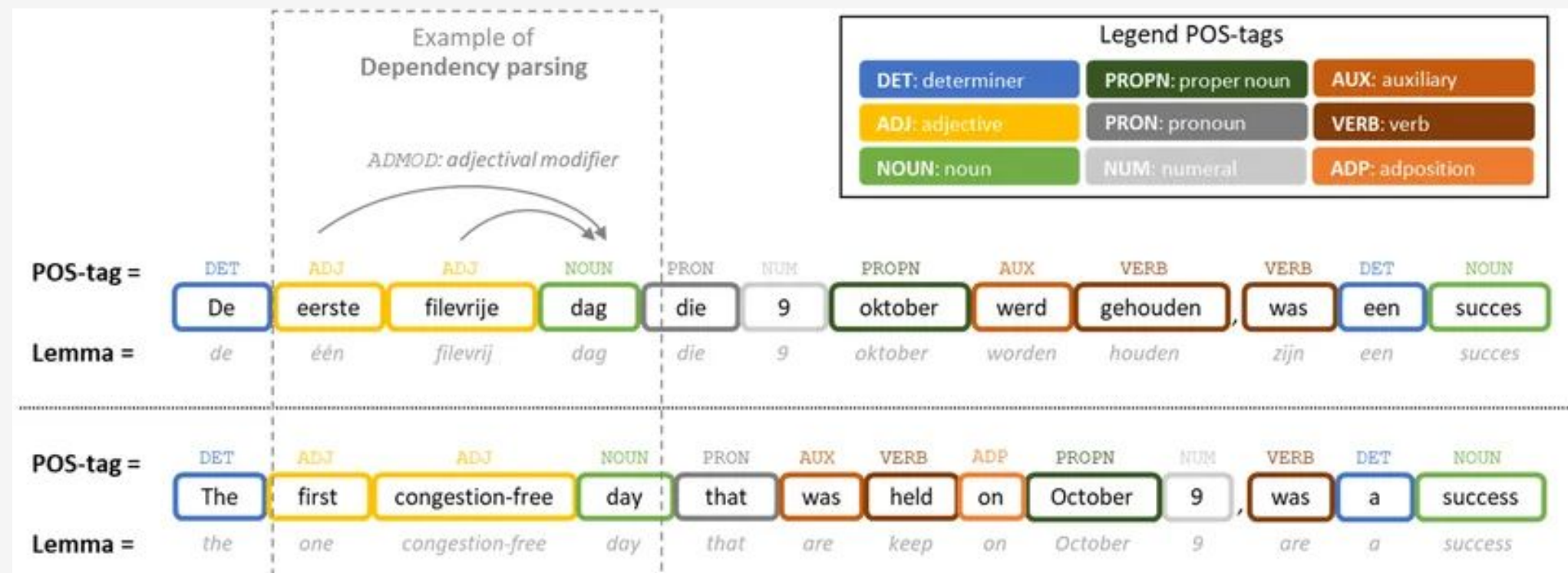


Задачи: морфология

Здесь несколько основных задач:

- ❖ POS-tagging;
- ❖ Лемматизация;
- ❖ Морфологический парсинг;
- ❖ Синтаксис (не морфология, но чаще всего вместе с ней).

К каким типам задач они относятся?



Задачи: семантика

- ❖ Word Sense Disambiguation / Induction - контекстное снятие омонимии (мы либо знаем о наборах значений у слов, либо нет);
- ❖ (Aspect-Based) Sentiment Analysis: определение эмоциональной окраски (относительно конкретного аспекта / в целом);
- ❖ Named Entity Recognition: выделение именованных сущностей;
- ❖ Natural Language Inference: определение связаны ли два текста (связаны / противоречат / нейтральны);
- ❖ Abstractive / Extractive Question Answering: ответ на вопрос (генеративно / через извлечение релевантного куска текста);
- ❖ Машинный перевод;
- ❖ Суммаризация;
- ❖ Разрешение анафоры / кореференции.

NER

Почему это сложно?

- ❖ Особенности орфографии: с большой буквы пишутся не только ИС;
- ❖ Многозначность: одна и та же ИС может иметь разный смысл;
- ❖ Необходимость учитывать контекст;
- ❖ Большой уровень разнообразия;
- ❖ Сложность получения чистых и разнообразных данных для обучения;
- ❖ Вложенность ИС как отдельная сложность и задача

B - beginning - 1 token сущности

I - inside - не 1/последний token сущности

O - out - не сущность

E - ending - последний token сущности

S - single - сущность из одного токена

Рауль	Модесто	Кастро	Рус	родился	в	городке	Биран	в	1931	году	.
B-PER	I-PER	I-PER	E-PER	OUT	OUT	OUT	S-LOC	OUT	OUT	B-DATE	E-DATE

Рауль Модесто Кастро Рус (Raul Modesto Castro Ruz) родился 3 июня 1931 года в городке

Биран в кубинской провинции Ориенте (нынешняя территория провинции Ольгин) , в семье

СПАНЫ

УПОМИНАНИЯ

✘ [Рауль][Модесто][Кастро]
[Рус]

Person

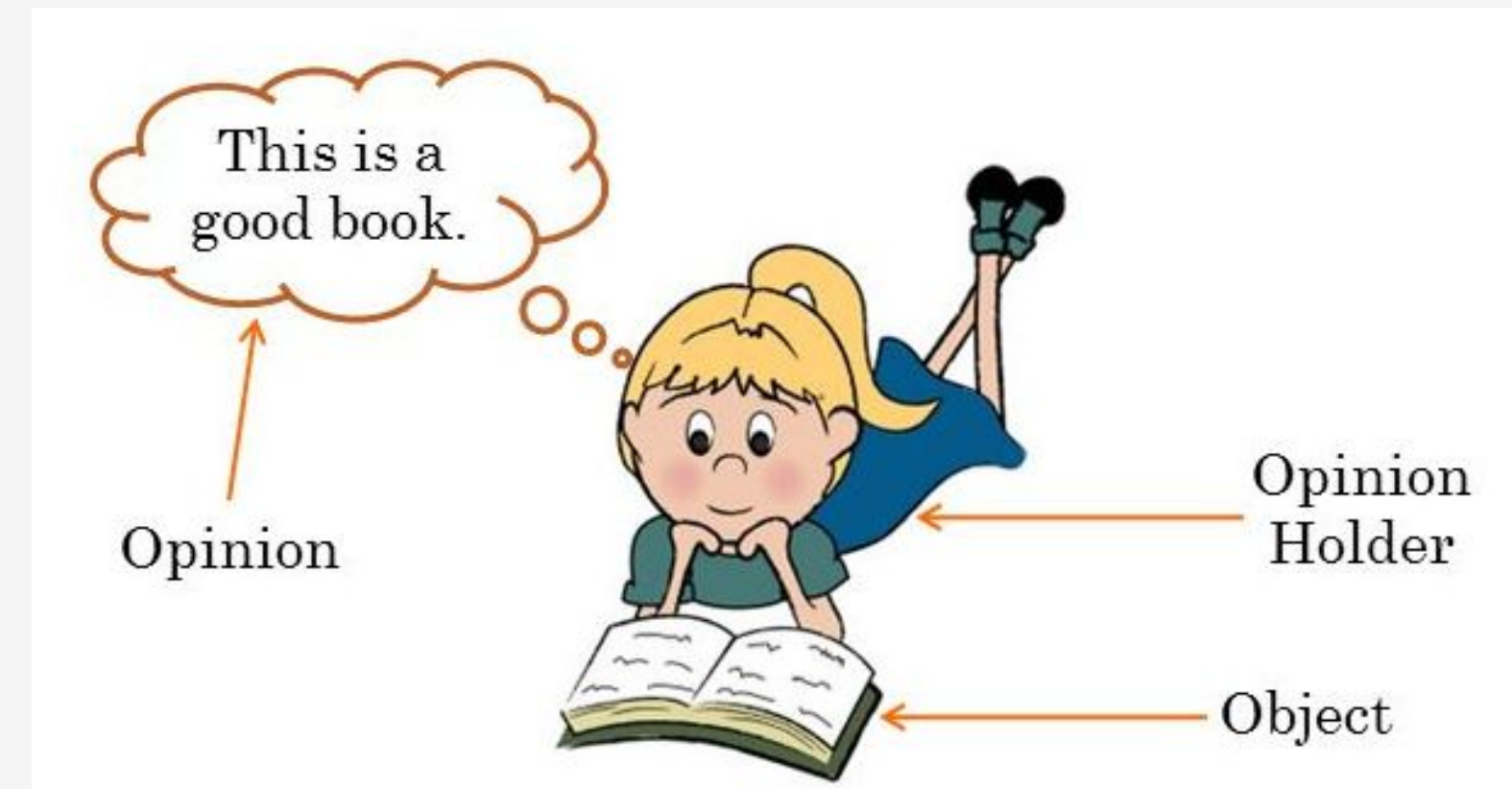
✘ [Raul][Modesto][Castro]
[Ruz]

Person

(AB)SA

Три близкие друг к другу задачи:

- ❖ Sentiment analysis: выявление / классификация оценочных суждений об объекте
- ❖ Opinion mining;
- ❖ Subjectivity detection: выявление субъективной информации (факты vs. мнения/эмоции).

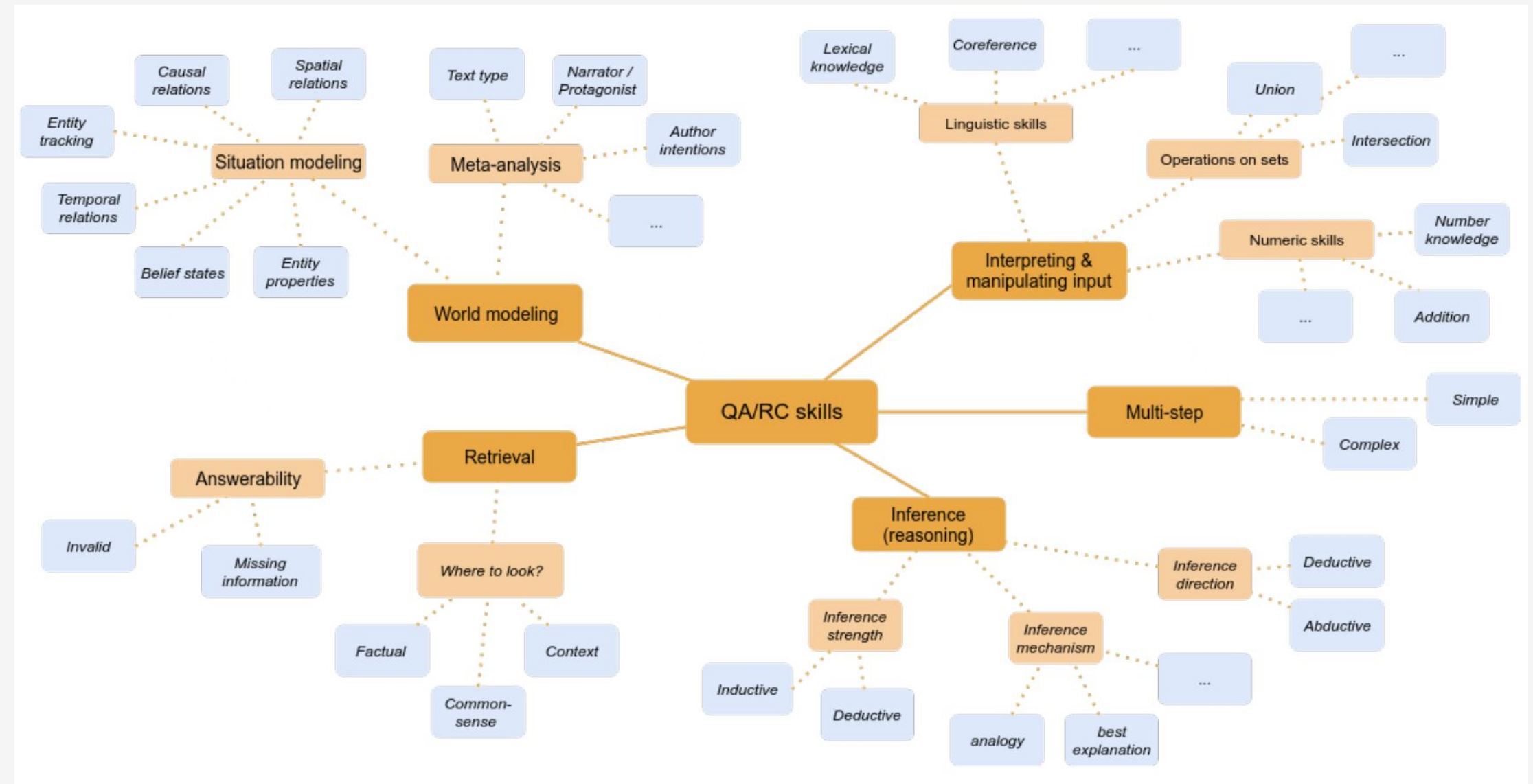


довольно_ADV приятный_ADJ впечатление_NOUN от_ADP Sushi_PROPN lounge_X _PUNCT это_PRON самый_ADJ дальний_ADJ
зал_NOUN _PUNCT заказывать_VERB с_ADP друг_NOUN по_ADP японский_ADJ меню_NOUN _PUNCT суша_NOUN
понравиться_VERB _PUNCT из_ADP основной_ADJ _PUNCT выбор_NOUN не_PART большой_ADJ _PUNCT зато_CCONJ
заказать_VERB из_ADP Италия_PROPN _PUNCT все_PART как_SCONJ обычно_ADV вкусный_ADJ _PUNCT _PUNCT
_PUNCT _PUNCT _PUNCT _PUNCT спасибо_NOUN так_ADV держать_VERB

QA

Почему это сложно:

- ❖ В вопросе могут присутствовать множественные сущности или связи;
- ❖ Время и другие обстоятельства могут быть не выражены эксплицитно;
- ❖ Могут требоваться дополнительные рассуждения по данным.



QA: метрики

Зачем свои метрики: правильный ответ не всегда дословно совпадает с эталоном.

Метрики:

- ❖ Полное совпадение
- ❖ BLEU, NIST (на основе BLEU, но каждой n-грамме присваивается вес в зависимости от встречаемости в языке), ROUGE (считаем не только precision, но и остальные две метрики), METEOR (учитываются не только точные совпадения слов, но и наличие однокоренных слов и синонимов);

$$\text{BLEU} = \min \left(1, \frac{\text{output-length}}{\text{reference-length}} \right) \left(\prod_{i=1}^4 \text{precision}_i \right)^{\frac{1}{4}}$$

- ❖ BERTScore: используем специального Берта, который обучен оценивать ответы (близко к NLI);
- ❖ LLM as a Judge: проверка сторонней, более "сильной" моделью (чаще всего LLM)
- ❖ Human feedback (на практике редко реализуемо).

Ключевые слова

Сложности:

- ❖ Мы знаем как минимум четыре способа понимать, что такое ключевые слова
- ❖ Разные люди по-разному выделяют слова даже в рамках одного определения
- ❖ Нужно понимание предметной области для разметки
- ❖ Может ли быть текст из одних ключевых слов? Вообще без ключевых слов? Какая плотность должна быть у разметки?

