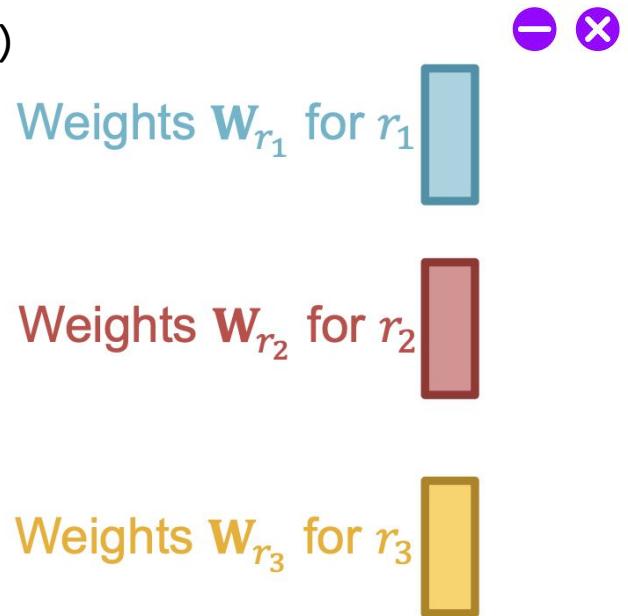
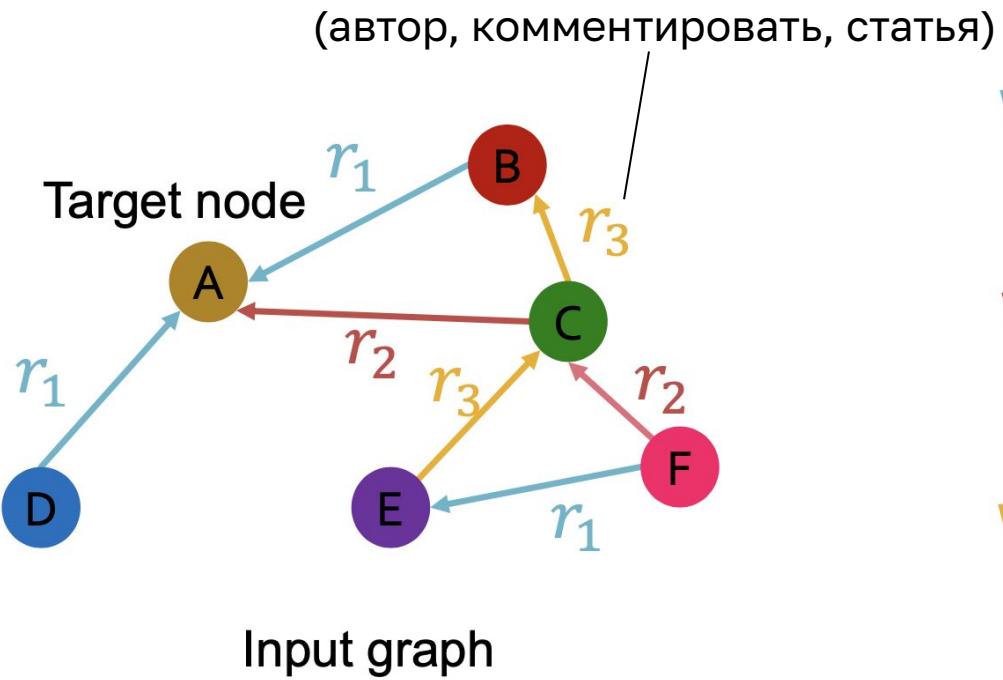


ИТМО

Анализ графовых данных и глубокое обучение

Азимов Рустам
Высшая школа цифровой культуры

В предыдущих сериях



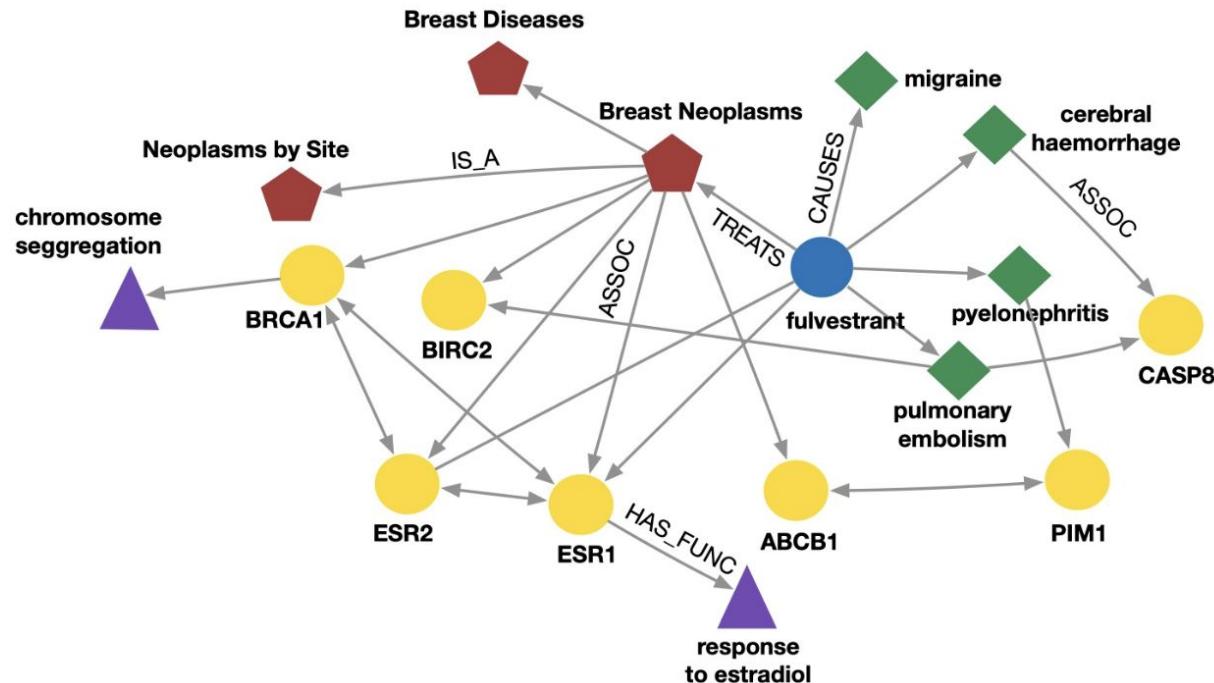
В этой лекции

- Графы знаний (**knowledge graphs**) – вид гетерогенных графов
- Вершины – сущности (**entities**) различных типов
- Рёбра описывают отношения (**relationships**) между сущностями
- Примеры графов знаний
 - Google Knowledge Graph
 - Yandex Object Answer
 - Facebook Graph API
- Примеры открытых графов знаний
 - Freebase, Wikidata, Dbpedia



Bio Knowledge Graph

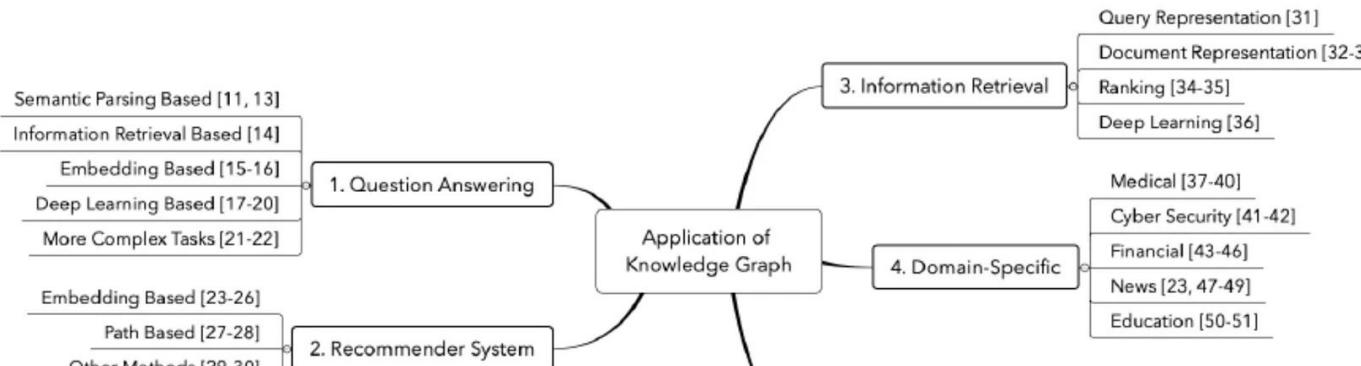
ViTMO



- Drug
- Disease
- Adverse event
- Protein
- Pathways

Применение графов знаний

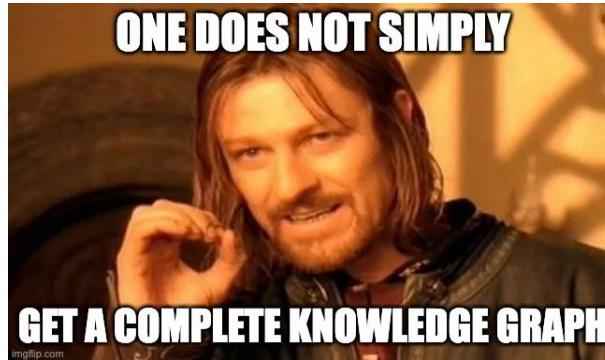
ІТМО



Характеристики графов знаний

ІТМО

- Графы знаний содержат миллионы вершин и рёбер
- Заполнять их всеми фактами очень трудоемко и в итоге получим полный граф, с которым тяжело работать
- Поэтому часто многие отношения упущены и очень актуальной становятся задачи **link prediction** и **knowledge graph completion**



Пример: Freebase

■ Freebase

- ~80 million **entities**
- ~38K **relation types**
- ~3 billion **facts/triples**



93.8% of persons from Freebase
have no place of birth and 78.5%
have no nationality!

■ Datasets: FB15k/FB15k-237

- A **complete** subset of Freebase, used by researchers to learn KG models

Dataset	Entities	Relations	Total Edges
FB15k	14,951	1,345	592,213
FB15k-237	14,505	237	310,079

[1] Paulheim, Heiko. "Knowledge graph refinement: A survey of approaches and evaluation methods." *Semantic web* 8.3 (2017): 489-508.

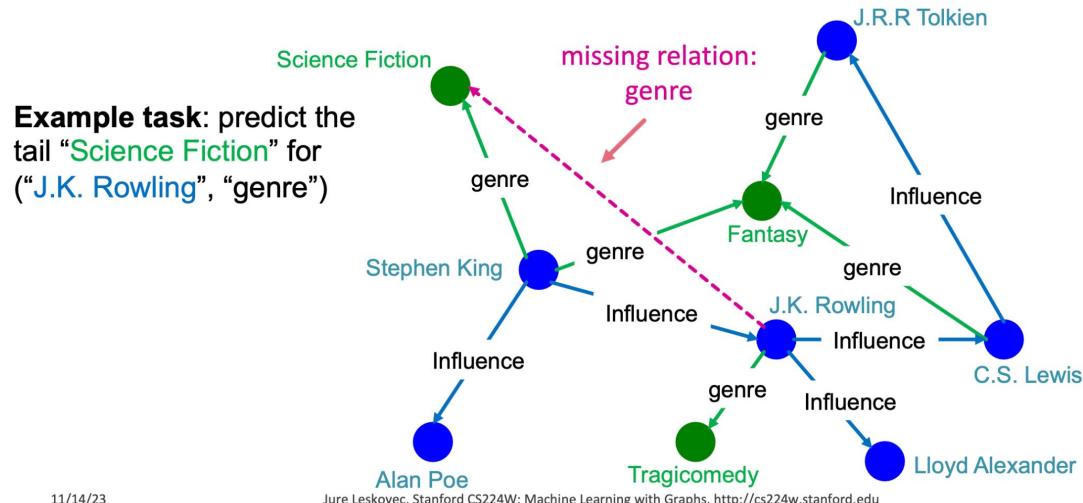
[2] Min, Bonan, et al. "Distant supervision for relation extraction with an incomplete knowledge base." *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*. 2013.

ViTMO

Knowledge Graph Completion

KG Completion

- Ребра в KG представлены в виде троек (h, r, t)
- Для заданных ($head, relation$) предсказать пропущенные tails



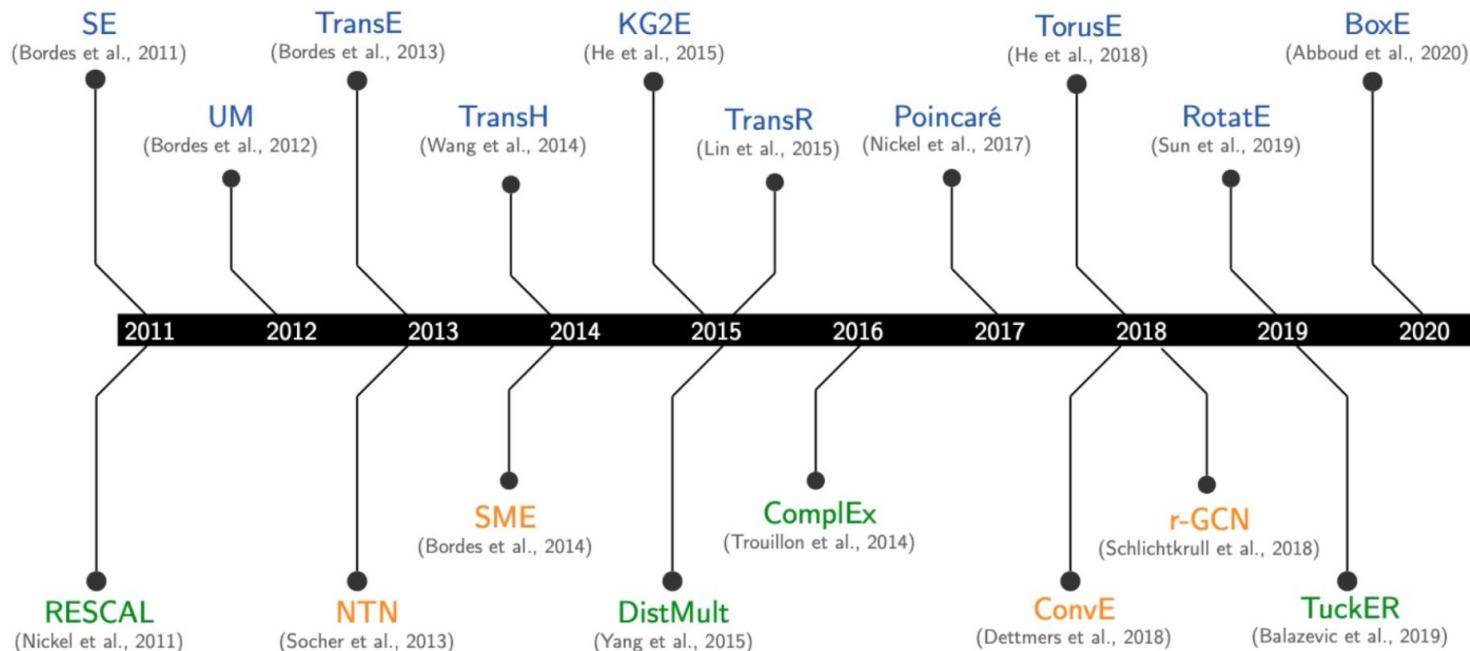
KG Completion

- Основная идея – использовать эмбеддинги из пространства небольшой размерности для представления сущностей и отношений
 - Можно использовать GNN, например RGCN
 - Сегодня речь пойдет о простом обучении эмбеддингов (без GNN)
- После обучения эмбеддингов хотим получить, что эмбеддинг (h, r) близок к эмбеддингу t
- Различные модели предлагают свои определения эмбеддингов (h, r) и функции ошибок



KG Embedding Models

VIITMO



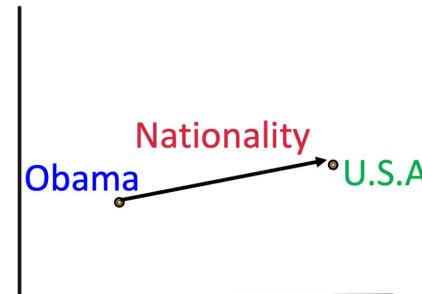
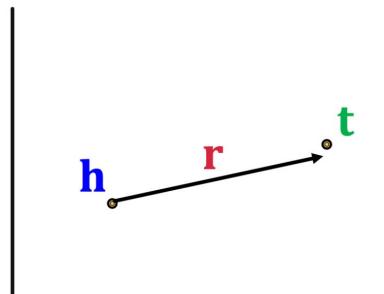
ViTMO

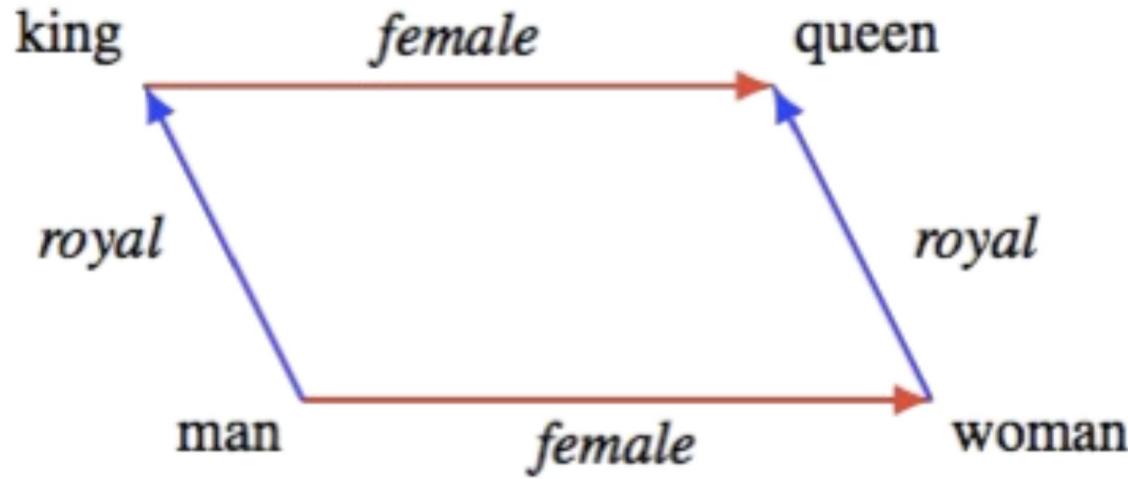
Knowledge Graph Completion: TransE

- Обучаем эмбеддинги для сущностей и отношений так, чтобы сумма эмбеддингов для h и r была близка к эмбеддингу t , если существует связь (h, r, t)

■ **TransE:** $h + r \approx t$ if the given link exists else $h + r \neq t$

Entity scoring function: $f_r(h, t) = -||h + r - t||$





TransE: Embedding Learning

Algorithm 1 Learning TransE

input Training set $S = \{(h, r, t)\}$, entities and rel. sets E and R , margin γ , embeddings dim. k .

1: **initialize** $r \leftarrow \text{uniform}(-\frac{6}{\sqrt{k}}, \frac{6}{\sqrt{k}})$ for each $r \in R$
2: $r \leftarrow r / \|r\|$ for each $r \in R$
3: $e \leftarrow \text{uniform}(-\frac{6}{\sqrt{k}}, \frac{6}{\sqrt{k}})$ for each entity $e \in E$

4: **loop**

5: $e \leftarrow e / \|e\|$ for each entity $e \in E$

6: $S_{batch} \leftarrow \text{sample}(S, b)$ // sample a minibatch of size b

7: $T_{batch} \leftarrow \emptyset$ // initialize the set of pairs of triplets

8: **for** $(h, r, t) \in S_{batch}$ **do**

9: $(h', r, t') \leftarrow \text{sample}(S'_{(h, r, t)})$ // sample a corrupted triplet

10: $T_{batch} \leftarrow T_{batch} \cup \{(h, r, t), (h', r, t')\}$

11: **end for**

12: Update embeddings w.r.t.

13: **end loop**

Initialize relations r and entities e uniformly, then normalize.
 γ is the margin.

Sample triplet (h', r, t) that does not appear in the KG.

d represents distance (negative of score)

$$\sum_{((h, r, t), (h', r, t')) \in T_{batch}} \nabla [\gamma + d(\mathbf{h} + \mathbf{r}, \mathbf{t}) - d(\mathbf{h}' + \mathbf{r}, \mathbf{t}')]_+$$

Contrastive loss: Favors lower distance (or higher score) for valid triplets, high distance (or lower score) for corrupted ones

Свойства отношений

- Отношения в KG могут обладать различными свойствами и связями друг с другом
 - Симметричность/Антисимметричность
 - Обратные отношения
 - Транзитивное отношение
 - Отношение один-ко-многим (1-to-N)
- Различные модели могут поддерживать или не поддерживать отношения такого вида
- Если в KG много отношений, которые не поддерживает выбранная модель, то лучше использовать другую

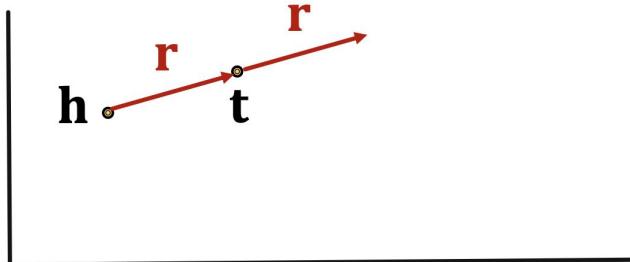


TransE: Антисимметричность

$$r(h, t) \Rightarrow \neg r(t, h) \quad \forall h, t$$



- **Example:** Hypernym (a word with a broader meaning: poodle vs. dog)
- **TransE** can model antisymmetric relations ✓
 - $h + r = t$, but $t + r \neq h$



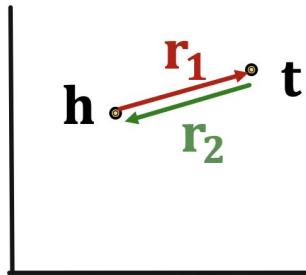
TransE: Обратные отношения

ІТМО

$$r_2(h, t) \Rightarrow r_1(t, h)$$



- Example : (Advisor, Advisee)
- TransE can model inverse relations ✓
- $\mathbf{h} + \mathbf{r}_2 = \mathbf{t}$, we can set $\mathbf{r}_1 = -\mathbf{r}_2$



TransE: Транзитивные отношения

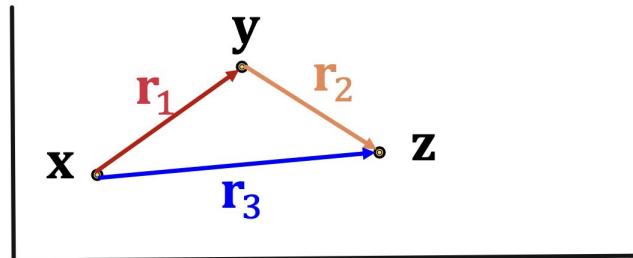
ІТМО

$$r_1(x, y) \wedge r_2(y, z) \Rightarrow r_3(x, z) \quad \forall x, y, z$$



- **Example:** My mother's husband is my father.
- **TransE** can model composition relations ✓

$$\mathbf{r}_3 = \mathbf{r}_1 + \mathbf{r}_2$$

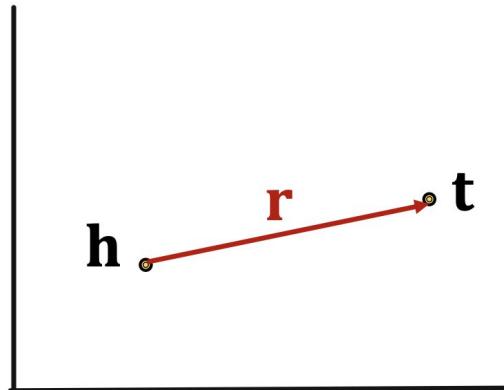


TransE: Симметричность

$$r(h, t) \Rightarrow r(t, h) \quad \forall h, t$$



- **Example:** Family, Roommate
- **TransE cannot** model symmetric relations ✗
only if $r = 0$, $h = t$

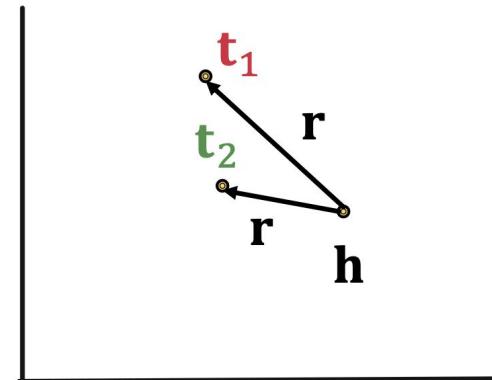


For all h, t that satisfy $r(h, t)$, $r(t, h)$ is also True, which means $\|h + r - t\| = 0$ and $\|t + r - h\| = 0$. Then $r = 0$ and $h = t$, however h and t are two different entities and should be mapped to different locations.

TransE: Отношения 1-to-N

- Example: (h, r, t_1) and (h, r, t_2) both exist in the knowledge graph, e.g., r is “StudentsOf” — ✗
- TransE cannot model 1-to-N relations ✗
 - t_1 and t_2 will map to the same vector, although they are different entities

- $t_1 = h + r = t_2$
- $t_1 \neq t_2$ contradictory!



TransE: Итоги



Model	Score	Embedding	Sym.	Antisym.	Inv.	Compos.	1-to-N
TransE	$-\ \mathbf{h} + \mathbf{r} - \mathbf{t}\ $	$\mathbf{h}, \mathbf{t}, \mathbf{r} \in \mathbb{R}^k$	✗	✓	✓	✓	✗

- Также существует модель RotatE (модификация TransE, использующая пространство комплексных чисел)

ViTMO

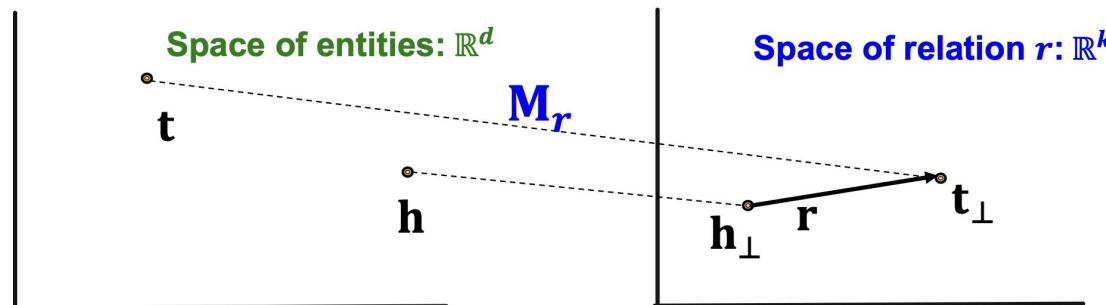
Knowledge Graph Completion: TransR

- Идея – использовать отдельное пространство эмбеддингов для отношений и проецировать в него эмбеддинги вершин умножением на матрицу M



- $h_{\perp} = M_r h, t_{\perp} = M_r t$
- **Score function:** $f_r(h, t) = -||h_{\perp} + r - t_{\perp}||$

Use M_r to project from entity space \mathbb{R}^d to relation space \mathbb{R}^k !



TransR: Симметричные отношения

ІТМО

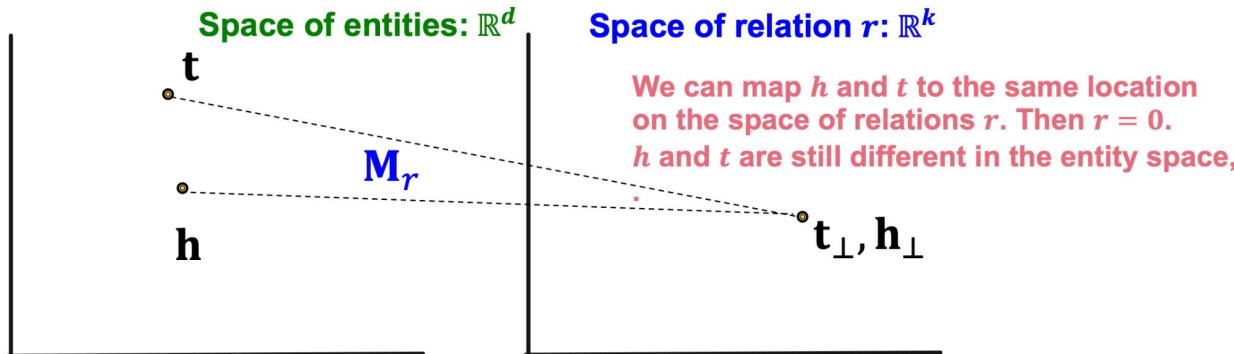
$$r(h, t) \Rightarrow r(t, h) \quad \forall h, t$$

Note different
symmetric
relations may
have different M_r



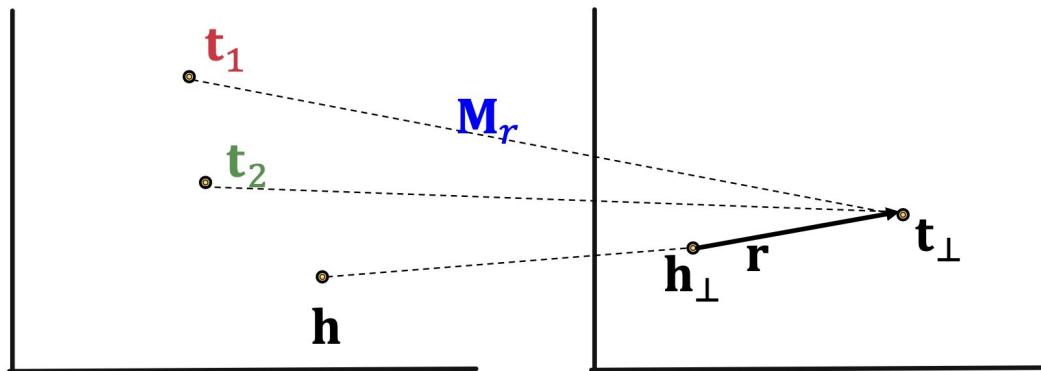
- **Example:** Family, Roommate
- TransR can model symmetric relations

$$\mathbf{r} = 0, \mathbf{h}_\perp = \mathbf{M}_r \mathbf{h} = \mathbf{M}_r \mathbf{t} = \mathbf{t}_\perp \checkmark$$



TransR: Отношения 1-to-N

- **Example:** If (h, r, t_1) and (h, r, t_2) exist in the knowledge graph.
- **TransR** can model 1-to-N relations ✓
 - We can learn \mathbf{M}_r so that $t_\perp = \mathbf{M}_r t_1 = \mathbf{M}_r t_2$
 - Note that t_1 does not need to be equal to t_2 !



TransR: Итоги



Model	Score	Embedding	Sym.	Antisym.	Inv.	Compos.	1-to-N
TransE	$-\ \mathbf{h} + \mathbf{r} - \mathbf{t}\ $	$\mathbf{h}, \mathbf{t}, \mathbf{r} \in \mathbb{R}^k$	✗	✓	✓	✓	✗
TransR	$-\ \mathbf{M}_r \mathbf{h} + \mathbf{r} - \mathbf{M}_r \mathbf{t}\ $	$\mathbf{h}, \mathbf{t} \in \mathbb{R}^k,$ $\mathbf{r} \in \mathbb{R}^d,$ $\mathbf{M}_r \in \mathbb{R}^{d \times k}$	✓	✓	✓	✓	✓

- Также может быть показано, что остальные виды отношений поддерживаются в TransR как и в TransE

ViTMO

Knowledge Graph Completion: DistMult

- TransE и TransR использовали минус L1/L2 норму в качестве метрики
- Попробуем усложнить

Score function: $f_r(h, t) = h \cdot A \cdot t$

$$h, t \in \mathbb{R}^k, A \in \mathbb{R}^{k \times k}$$

- Слишком сложно, легко переобучается
 - Упростить, сделав A – диагональной матрицей (DistMult)



- **Score function:** $f_r(h, t) = \langle h, r, t \rangle = \sum_i h_i \cdot r_i \cdot t_i$



- $h, r, t \in \mathbb{R}^k$

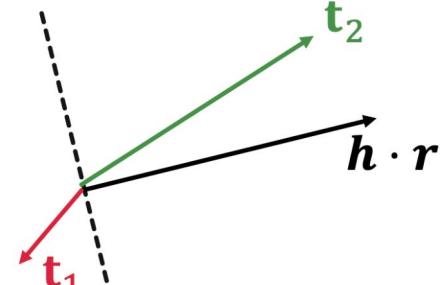
- **Intuition of the score function:** Can be viewed as a **cosine similarity** between $h \cdot r$ and t

where $h \cdot r$ is defined as $[h \cdot r]_i = h_i \cdot r_i$

- **Example:**

Hadamard product

$$f_r(h, t_1) < 0, \quad f_r(h, t_2) > 0$$



DistMult: Итоги



Model	Score	Embedding	Sym.	Antisym.	Inv.	Compos.	1-to-N
TransE	$-\ \mathbf{h} + \mathbf{r} - \mathbf{t}\ $	$\mathbf{h}, \mathbf{t}, \mathbf{r} \in \mathbb{R}^k$	✗	✓	✓	✓	✗
TransR	$-\ M_r \mathbf{h} + \mathbf{r} - M_r \mathbf{t}\ $	$\mathbf{h}, \mathbf{t} \in \mathbb{R}^k,$ $\mathbf{r} \in \mathbb{R}^d,$ $M_r \in \mathbb{R}^{d \times k}$	✓	✓	✓	✓	✓
DistMult	$\langle \mathbf{h}, \mathbf{r}, \mathbf{t} \rangle$	$\mathbf{h}, \mathbf{t}, \mathbf{r} \in \mathbb{R}^k$	✓	✗	✗	✗	✓

- Из-за коммутативности используемых в DistMult операций не поддерживаются антисимметричные, обратные и транзитивные отношения

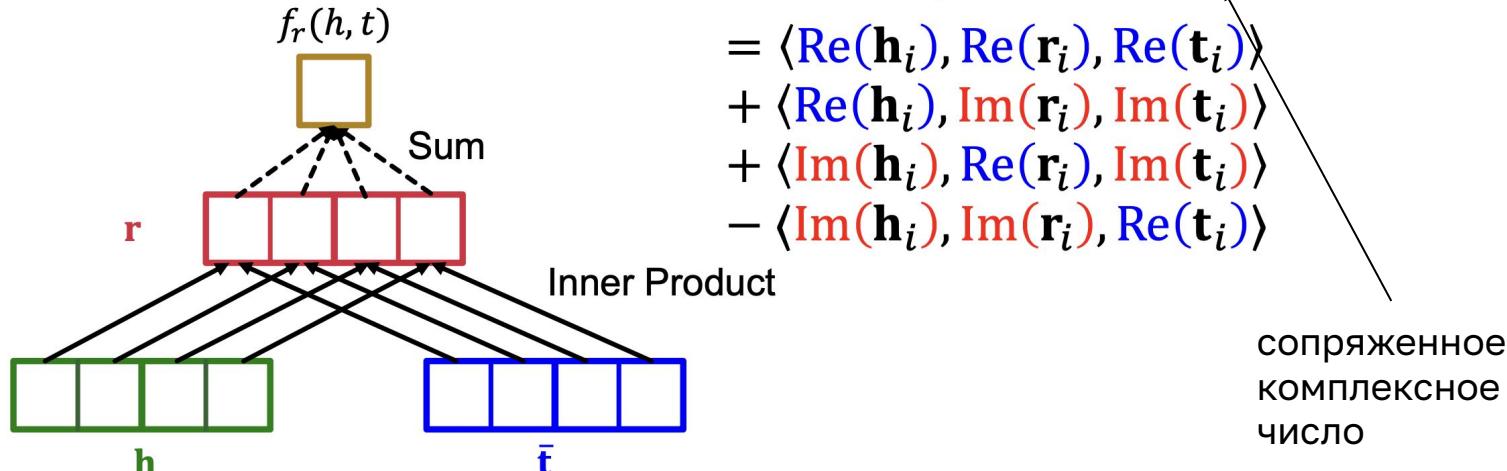
ViTMO

Knowledge Graph Completion: ComplEx

- ComplEx основан на DistMult, но использует комплексные числа и различает head и tail



■ **Score function** $f_r(h, t) = \operatorname{Re}(\sum_i \mathbf{h}_i \cdot \mathbf{r}_i \cdot \bar{\mathbf{t}}_i)$



KG Completion: Итоги

- При выборе модели учитывайте какие типы отношений важны для ваших данных



Model	Score	Embedding	Sym.	Antisym.	Inv.	Compos.	1-to-N
TransE	$-\ h + r - t\ $	$h, t, r \in \mathbb{R}^k$	✗	✓	✓	✓	✗
TransR	$-\ M_r h + r - M_r t\ $	$h, t \in \mathbb{R}^k,$ $r \in \mathbb{R}^d,$ $M_r \in \mathbb{R}^{d \times k}$	✓	✓	✓	✓	✓
DistMult	$\langle h, r, t \rangle$	$h, t, r \in \mathbb{R}^k$	✓	✗	✗	✗	✓
ComplEx	$\text{Re}(\langle h, r, \bar{t} \rangle)$	$h, t, r \in \mathbb{C}^k$	✓	✓	✓	✗	✓

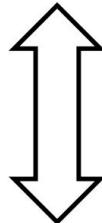
ViTMO

Reasoning over Knowledge Graphs

One-hop Query



- **KG completion:** Is link (h, r, t) in the KG?



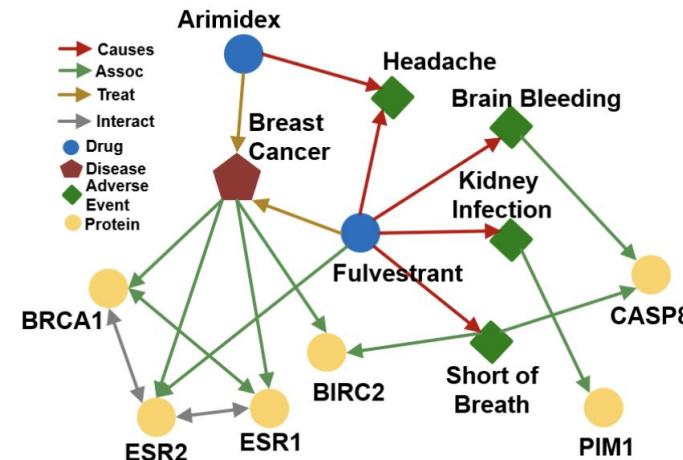
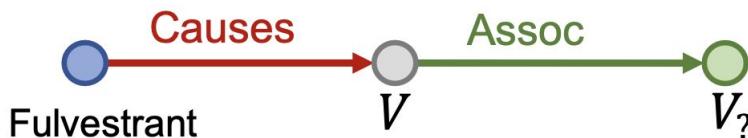
- **One-hop query:** Is t an answer to query (h, r) ?

Path Queries

- Рассмотрим более сложные запросы к графам знаний



Query: (e:Fulvestrant, (r:Causes, r:Assoc))

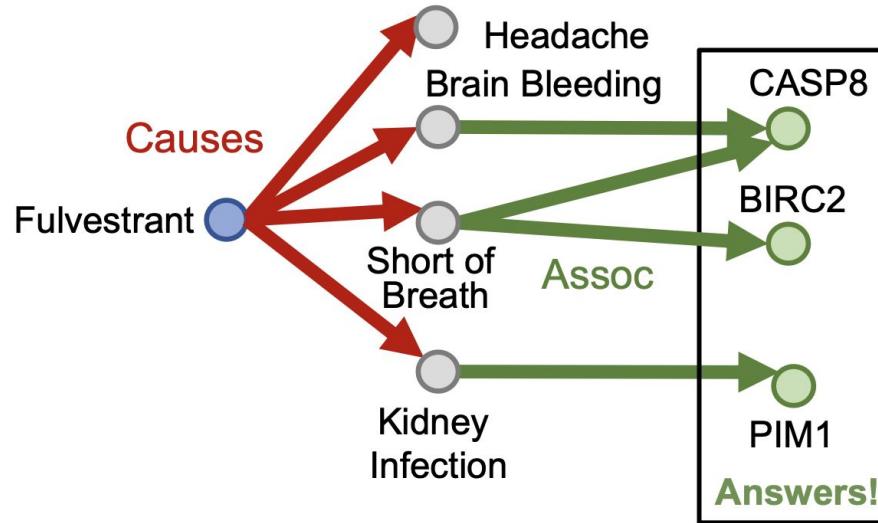


Path Queries Evaluation

- Можно отвечать на запрос, обходя граф



Query: (e:Fulvestrant, (r:Causes, r:Assoc))



Path Queries Evaluation

- Но KGs содержат множество пропущенных связей
- Поэтому просто обходя граф мы упустим много вершин, которые должны быть в результирующем множестве
- Можно ли просто решить задачу KG Completion, а потом обходить граф и отвечать на path queries?
- Нет, KG станет слишком плотным графом для эффективной работы с ним
- Что делать?

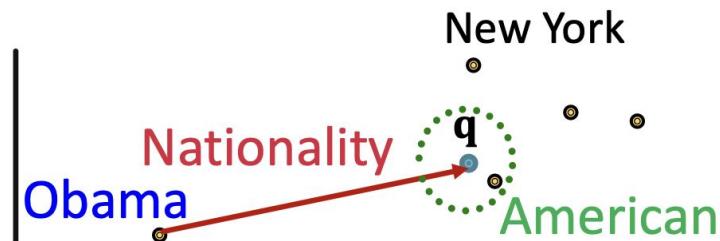
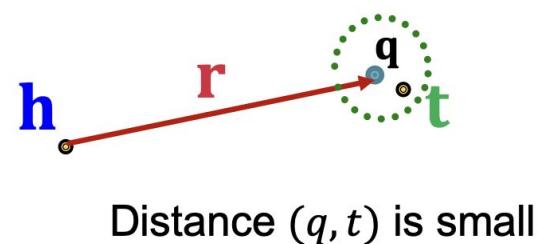


Predictive Queries

- Нужно отвечать на запросы, учитывая возможное отсутствие связей в графе
- Path queries
 - Обобщаем TransE
- Conjunctive queries
 - Query2box
- Boolean (and + or) queries
 - Приведение к нормальной форме, затем Query2box



- Запрос (эмбеддинг) $q = h + r$, должен быть близок к ответам t (эмбеддингам)



Path Queries: TransE

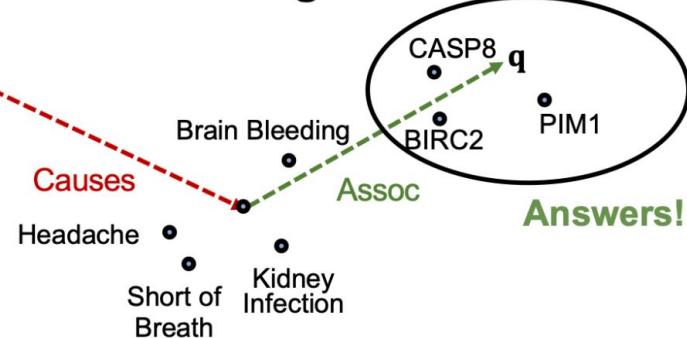
- TransE может естественным образом отвечать на path queries, так как поддерживает транзитивные отношения
- В отличие от DistMult/ComplEx



Query Plan



Embedding Process



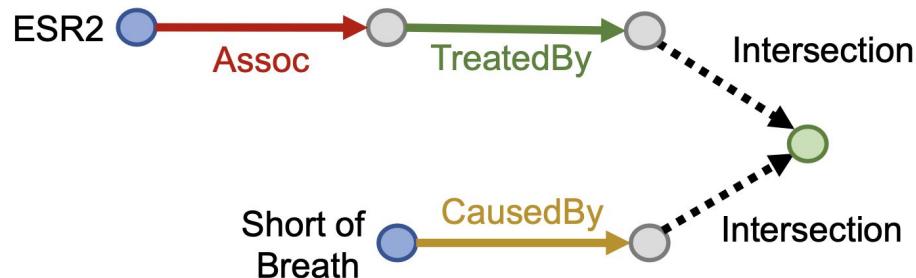
Answers!

Conjunctive Queries

- Более сложные запросы с конъюнкцией
- На каждом шаге хотим выявлять множество сущностей, близких к текущему подзапросу
- В момент применения конъюнкции хотим как то пересекать эти множества



Query plan:



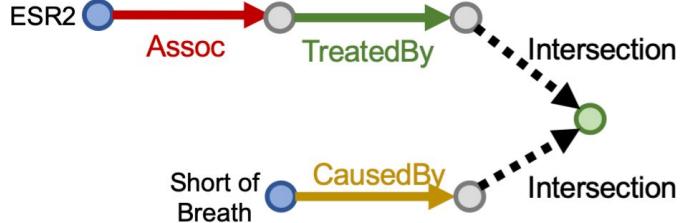
Conjunctive Queries: Query2box

ІТМО

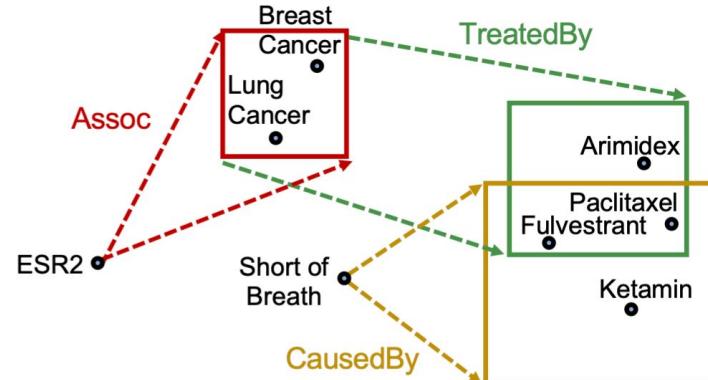
- Каждому запросу ставим в соответствие гиперпрямоугольник
- Эмбеддинги сущностей – точки, а отношения переводят одни прямоугольники в другие (projection operator)



Query Plan



Embedding Space



Query2box: Projection Operator

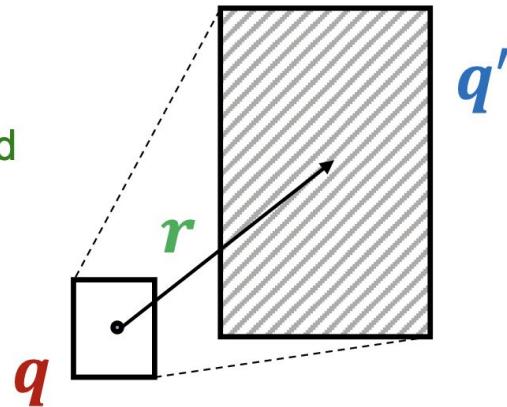
- $\mathcal{P} : \text{Box} \times \text{Relation} \rightarrow \text{Box}$

- ×

$$\text{Cen}(q') = \text{Cen}(q) + \text{Cen}(r)$$

$$\text{Off}(q') = \text{Off}(q) + \text{Off}(r)$$

"x" (cross) means the projection operator is a relation from any box and relation to a new box



Query2box: Intersection Operator

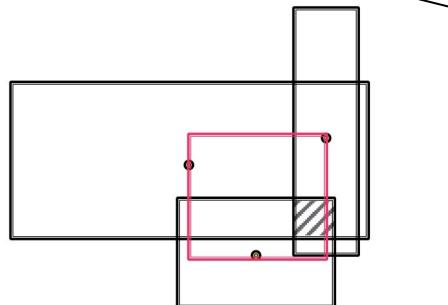
ИТМО

- $\mathcal{I} : \text{Box} \times \dots \times \text{Box} \rightarrow \text{Box}$

Hadamard product
(element-wise product)

$$Cen(q_{inter}) = \sum_i \mathbf{w}_i \odot Cen(q_i)$$

$$\mathbf{w}_i = \frac{\exp(f_{cen}(Cen(q_i)))}{\sum_j \exp(f_{cen}(Cen(q_j)))} \quad Cen(q_i) \in \mathbb{R}^d$$
$$\mathbf{w}_i \in \mathbb{R}^d$$



Query2box: Intersection Operator

ІТМО



- $\mathcal{I} : \text{Box} \times \dots \times \text{Box} \rightarrow \text{Box}$

$$Off(q_{inter})$$

$$= \min(Off(q_1), \dots, Off(q_n))$$

$$\odot \sigma(f_{off}(Off(q_1), \dots, Off(q_n)))$$

guarantees shrinking

нейронная сеть

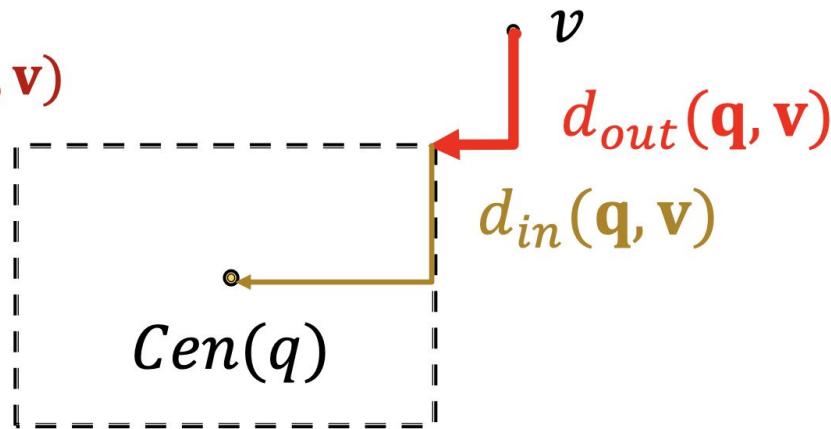
Query2box: Score Function

$$d_{box}(\mathbf{q}, \mathbf{v}) = d_{out}(\mathbf{q}, \mathbf{v}) + \alpha \cdot d_{in}(\mathbf{q}, \mathbf{v})$$



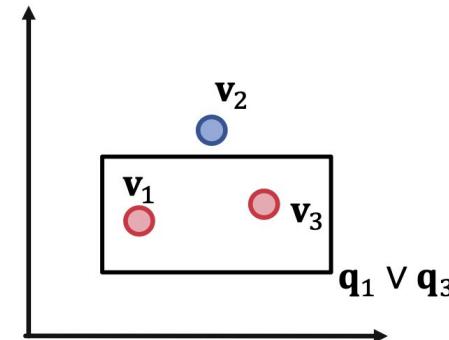
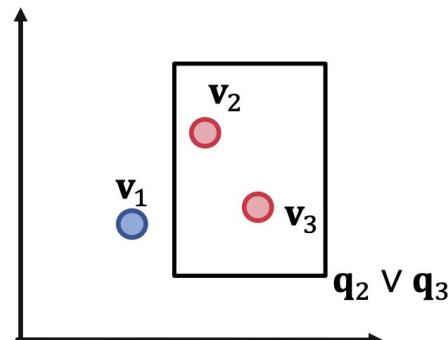
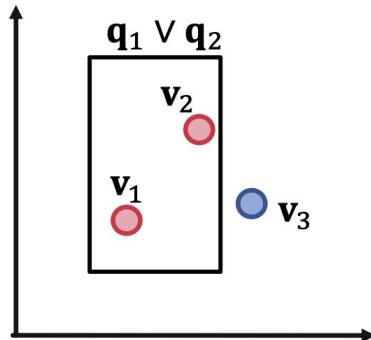
where $0 < \alpha < 1$.

$$f_q(v) = -d_{box}(\mathbf{q}, \mathbf{v})$$



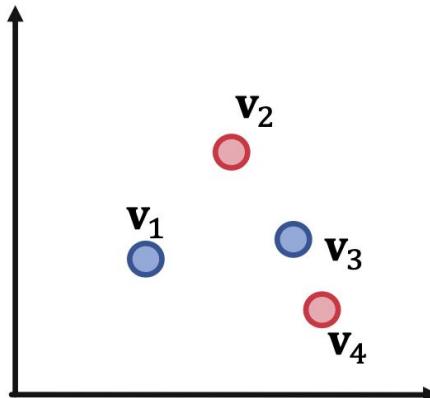
Boolean Queries

- Добавим возможность использовать дизъюнкцию в запросах



Boolean Queries

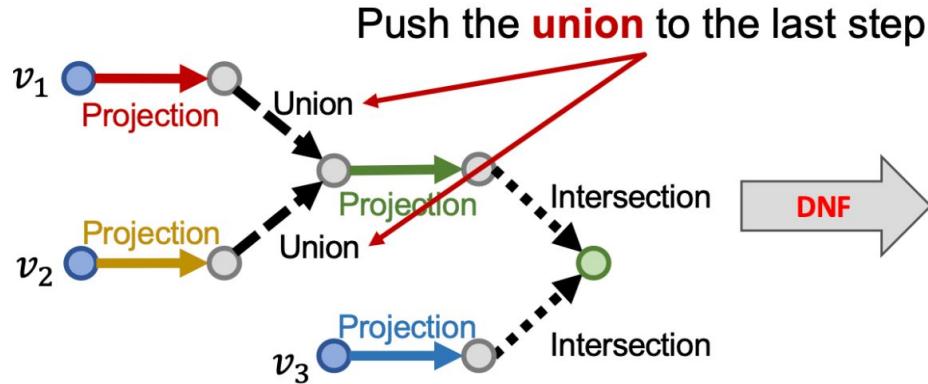
- Для того, чтобы различать много произвольных Булевых запросов
пространства эмбеддингов малой размерности может не хватить



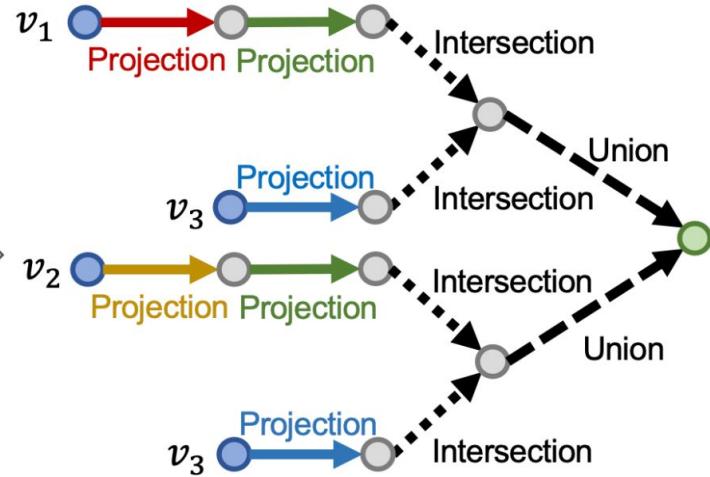
Query2box: Disjunctive Normal Form ViTMO



Original Query Plan



Converted Query Plan



Query2box: Disjunctive Normal Form

$$d_{box}(\mathbf{q}, \mathbf{v}) = \min(d_{box}(\mathbf{q}_1, \mathbf{v}), \dots, d_{box}(\mathbf{q}_m, \mathbf{v}))$$


- The process of embedding any AND-OR query q
 1. Transform q to equivalent DNF $q_1 \vee \dots \vee q_m$
 2. Embed q_1 to q_m
 3. Calculate the (box) distance $d_{box}(\mathbf{q}_i, \mathbf{v})$
 4. Take the minimum of all distance
 5. The final score $f_q(v) = -d_{box}(\mathbf{q}, \mathbf{v})$

Query2box: Training

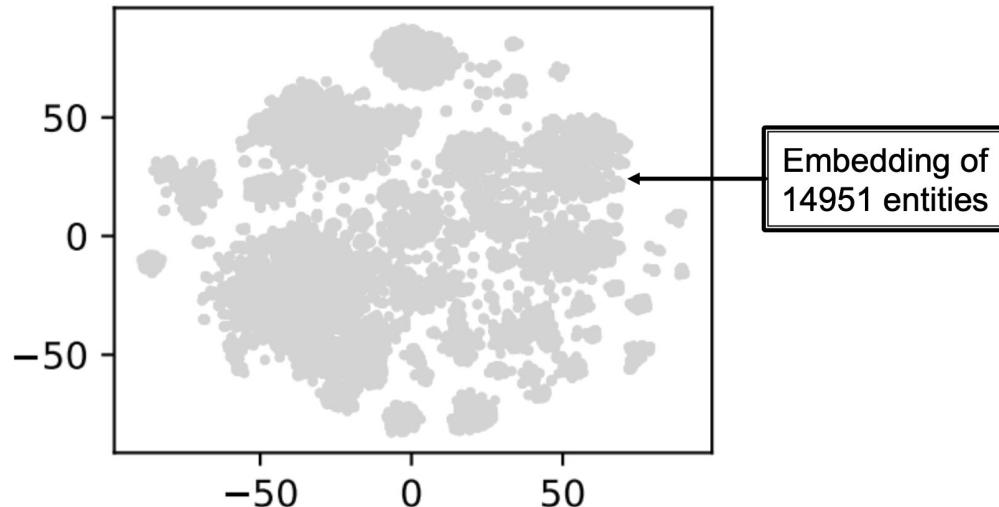
Training:



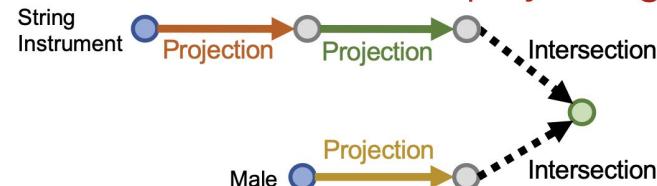
1. Sample a query q from the training graph G_{train} ,
answer $v \in \llbracket q \rrbracket_{G_{train}}$, and non-answer $v' \notin \llbracket q \rrbracket_{G_{train}}$
2. Embed the query q .
 - Use current operators, to compute query embedding.
3. Calculate the score $f_q(v)$ and $f_q(v')$.
4. Optimize embeddings and operators to minimize the
loss ℓ (maximize $f_q(v)$ while minimize $f_q(v')$):

$$\ell = -\log \sigma(f_q(v)) - \log(1 - \sigma(f_q(v')))$$

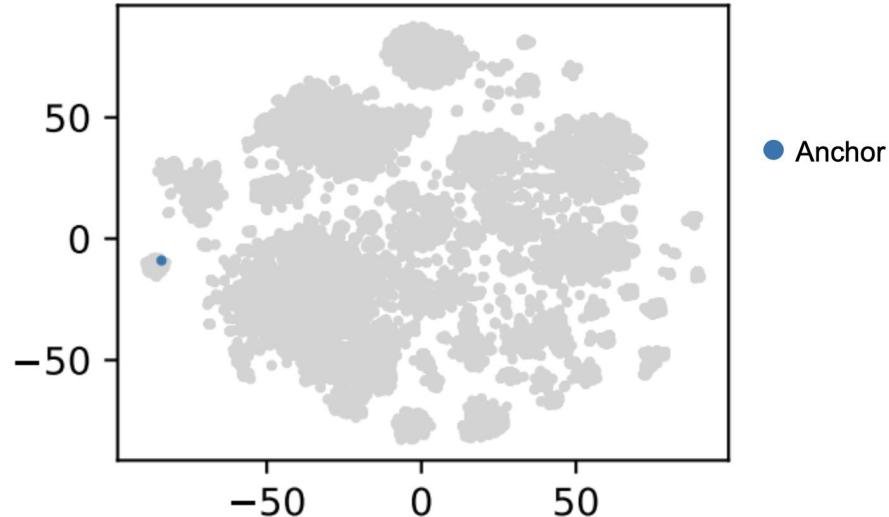
Query2box: Example



“List male instrumentalists who play string instruments”



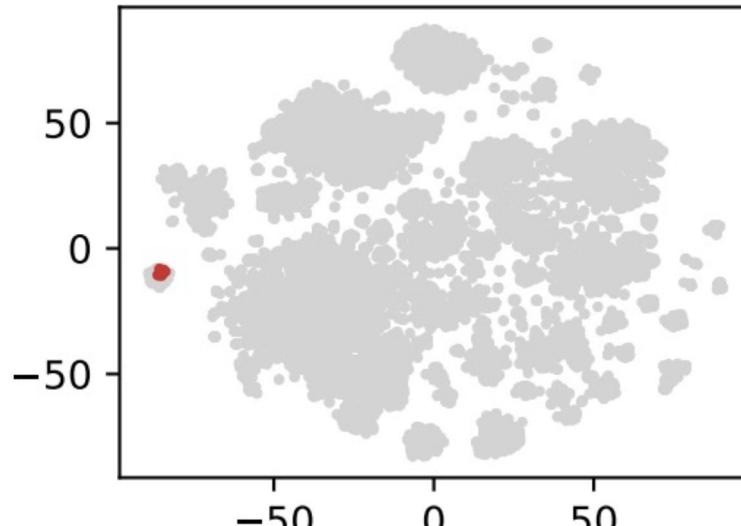
Query2box: Example



“List male instrumentalists who play string instruments”

String
Instrument

Query2box: Example



of string instruments: 10



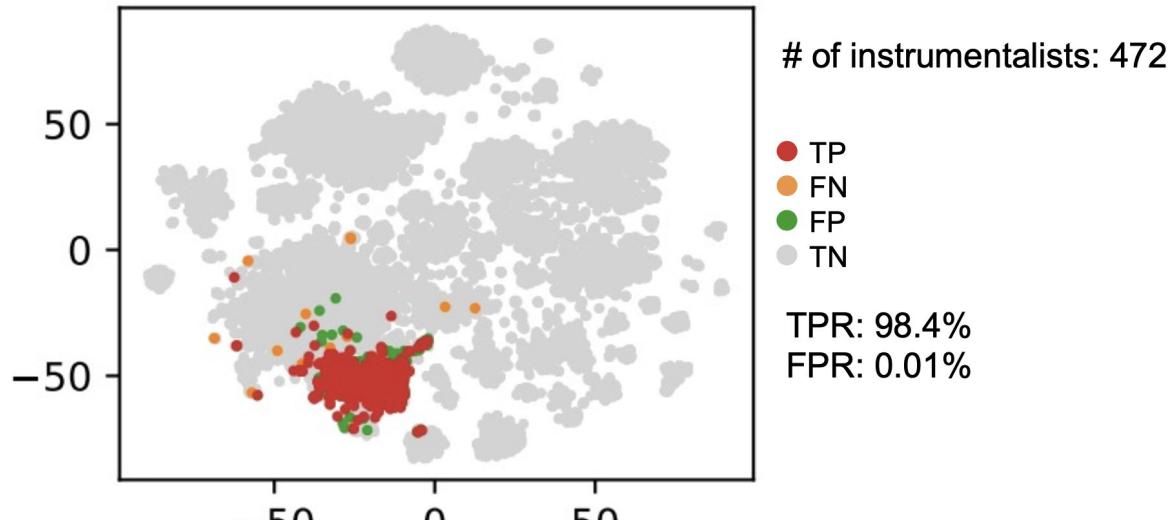
- TP
- FN
- FP
- TN

TPR: 100%
FPR: 0%

“List male instrumentalists who play string instruments”



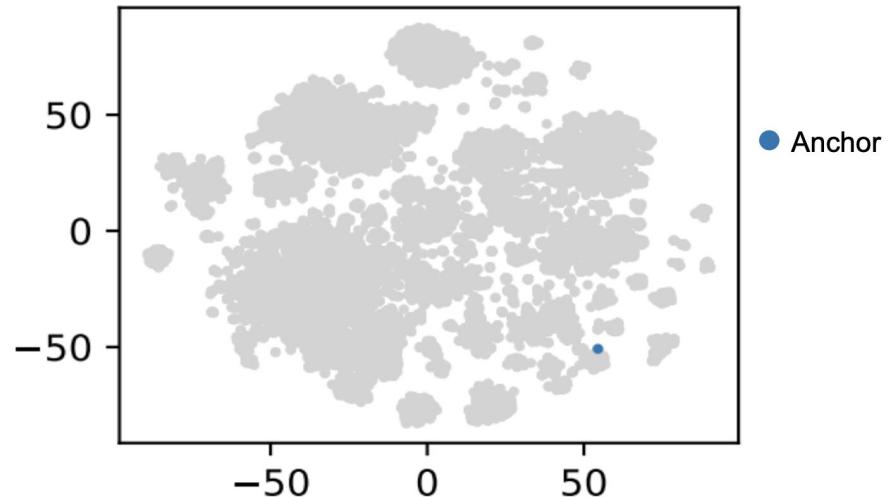
Query2box: Example



“List male instrumentalists who play string instruments”



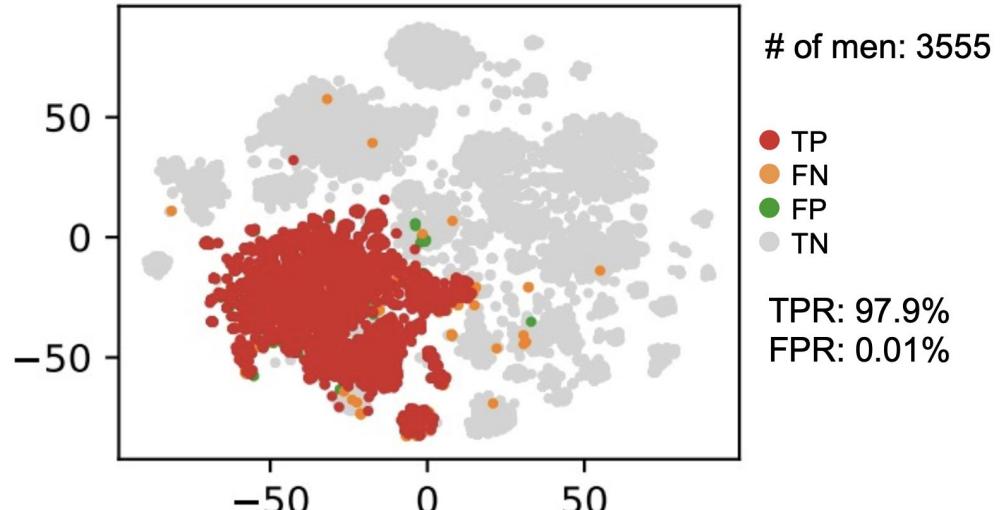
Query2box: Example



“List male instrumentalists who play string instruments”

Male

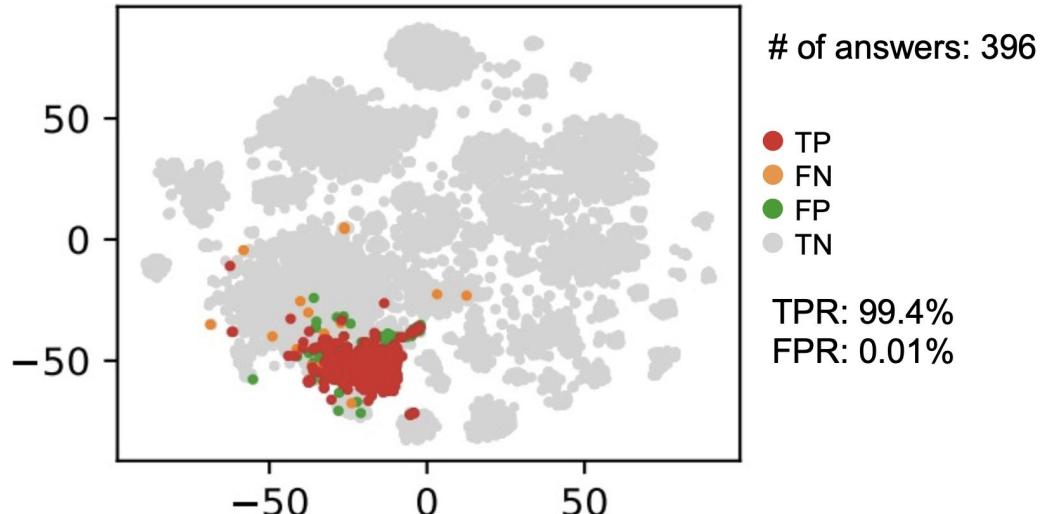
Query2box: Example



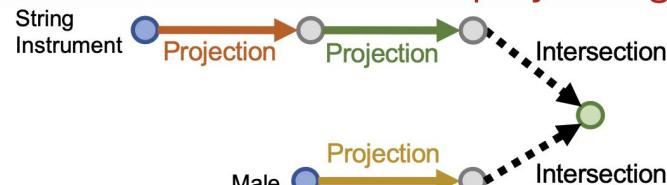
“List male instrumentalists who play string instruments”



Query2box: Example



"List male instrumentalists who play string instruments"



Заключение

Model	Score	Embedding	Sym.	Antisym.	Inv.	Compos.	1-to-N
TransE	$-\ \mathbf{h} + \mathbf{r} - \mathbf{t}\ $	$\mathbf{h}, \mathbf{t}, \mathbf{r} \in \mathbb{R}^k$	✗	✓	✓	✓	✗
TransR	$-\ M_r \mathbf{h} + \mathbf{r} - M_r \mathbf{t}\ $	$\mathbf{h}, \mathbf{t} \in \mathbb{R}^k$, $\mathbf{r} \in \mathbb{R}^d$, $M_r \in \mathbb{R}^{d \times k}$	✓	✓	✓	✓	✓
DistMult	$\langle \mathbf{h}, \mathbf{r}, \mathbf{t} \rangle$	$\mathbf{h}, \mathbf{t}, \mathbf{r} \in \mathbb{R}^k$	✓	✗	✗	✗	✓
ComplEx	$\text{Re}(\langle \mathbf{h}, \mathbf{r}, \bar{\mathbf{t}} \rangle)$	$\mathbf{h}, \mathbf{t}, \mathbf{r} \in \mathbb{C}^k$	✓	✓	✓	✗	✓



- Нужно отвечать на запросы, учитывая возможное отсутствие связей в графе
- Path queries – обобщаем TransE
- Conjunctive queries – Query2box
- Boolean (and + or) queries – DNF + Query2box