

# Модуль оценки качества векторного представления графов

Сравнение с альтернативными реализациями:

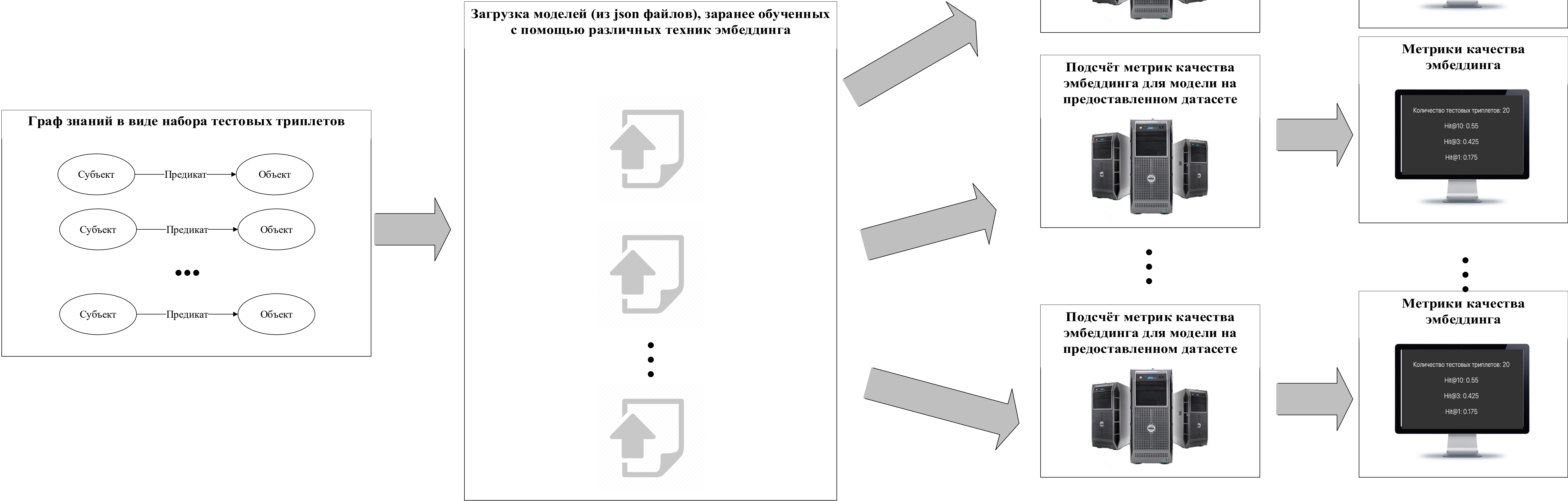
Параметр	OpenKE	Mana-ysh	Модуль KGE	$\alpha$
Реализация базовых техник эмбединга	1	0,6	1	0,25
Возможность получения метрики среднего ранга	1	0	1	0,2
Возможность получения метрик типа Hit@N	0	0	1	0,3
Наличие графического интерфейса	0	0	0,8	0,15
Возможность использования без установки	0	0	1	0,05
Качество документации	0,5	0,1	1	0,05
Итого	0,475	0,155	0,97	1

Цель работы:

разработать и предоставить пользователю модуль, позволяющий проводить оценку качества эмбединга (встраивания в векторное пространство) графа знаний. Модуль выдаёт пользователю метрики качества эмбединга графа на тестовом датасете для сравнения техник эмбединга и упрощения выбора подходящей техники. Также разработанный модуль с веб-интерфейсом можно использовать в качестве обучающего стенда при изучении эмбединга графа знаний.

Список решённых в работе задач:

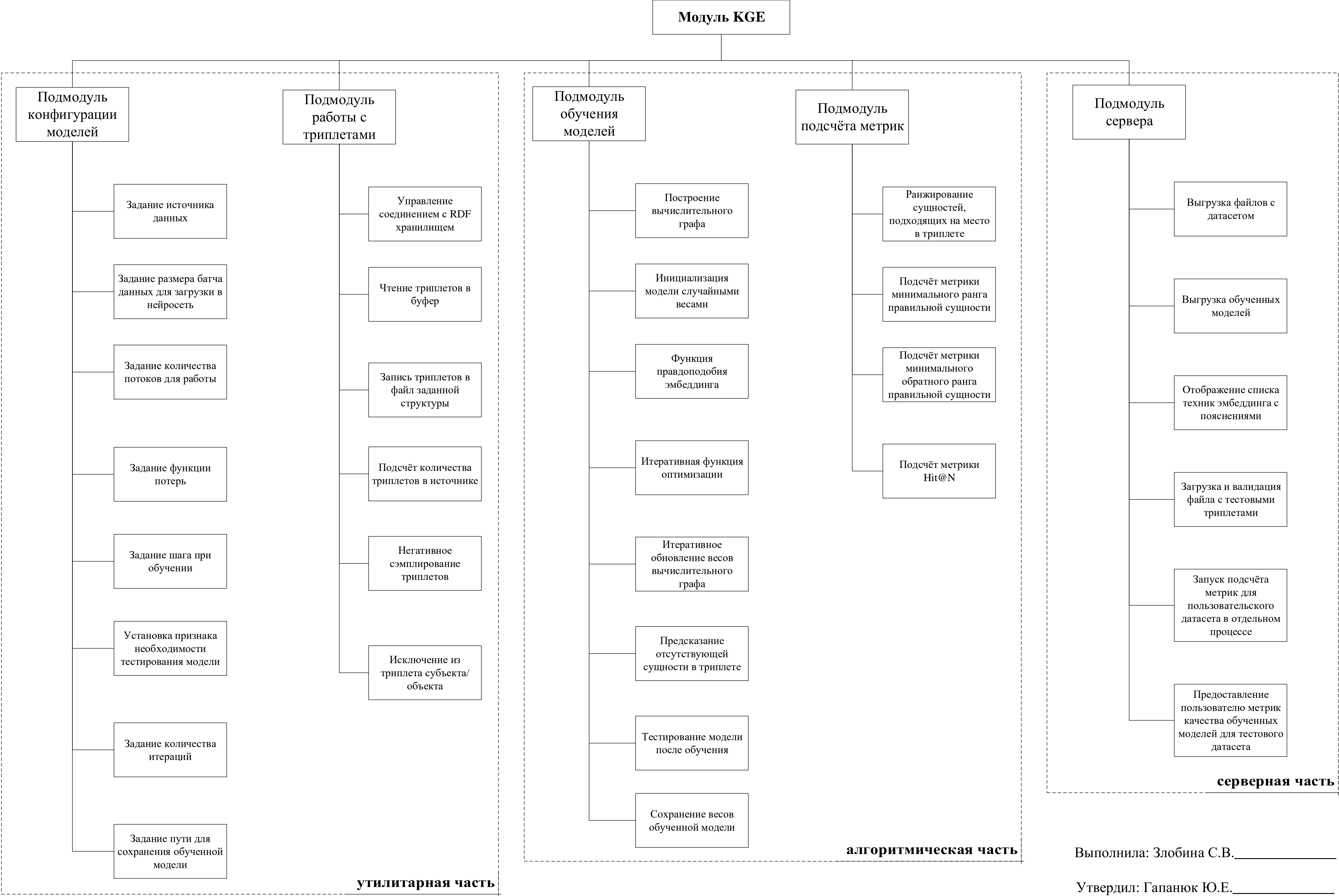
- Исследование предметной области и существующих техник эмбединга
- Выбор метрик и разработка алгоритмов для оценки качества эмбединга
- Объединение техник эмбединга и алгоритмов подсчёта метрик в общий модуль
- Разработка веб-интерфейса модуля
- Тестирование корректности работы модуля на различных датасетах



Выполнила: Злобина С.В.\_\_\_\_\_

Утвердил: Гапанюк Ю.Е.\_\_\_\_\_

# Структура модуля





# Эмбеддинг графа знаний

## Процесс эмбеддинга

### Эмбеддинг графа знаний -

это математическое преобразование графа знаний в вектор или в набор векторов в заданном векторном пространстве. Такое преобразование должно адекватно передавать семантику и топологию исходного графа.

1

**Выбор способа представления сущностей и их отношений в векторном пространстве.**

Сущности обычно представляются в виде векторов (точек в векторном пространстве), а отношения – в виде операторов в этом векторном пространстве.

2

**Задание функции правдоподобия преобразования (англ. scoring function).**

Для каждого триплета  $\langle h, r, t \rangle$  задаётся функция правдоподобия  $f_r(h, t)$ , чем больше значение этой функции, тем более вероятно то, что факт, описываемый триплетом, является истинным.

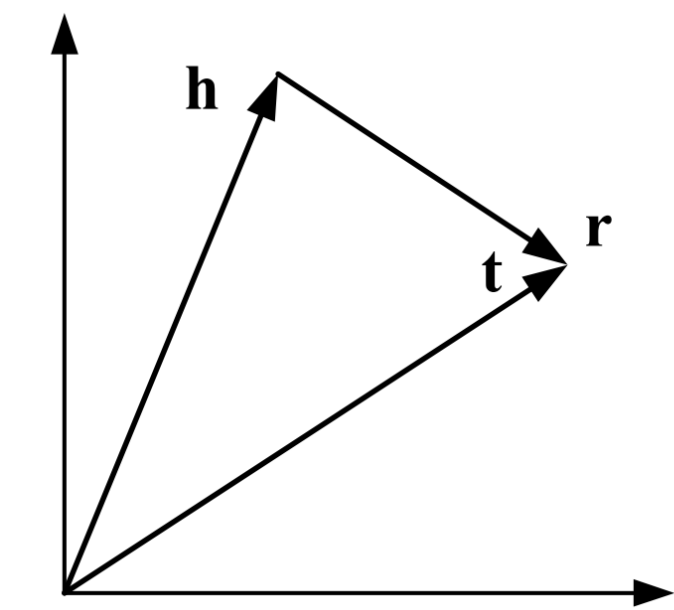
3

**Обучение модели эмбеддинга.**

Подразумевает решение оптимизационной задачи, заключающейся в максимизации суммарной функции правдоподобия для всех триплетов, содержащихся в графе знаний.

## Техники эмбеддинга

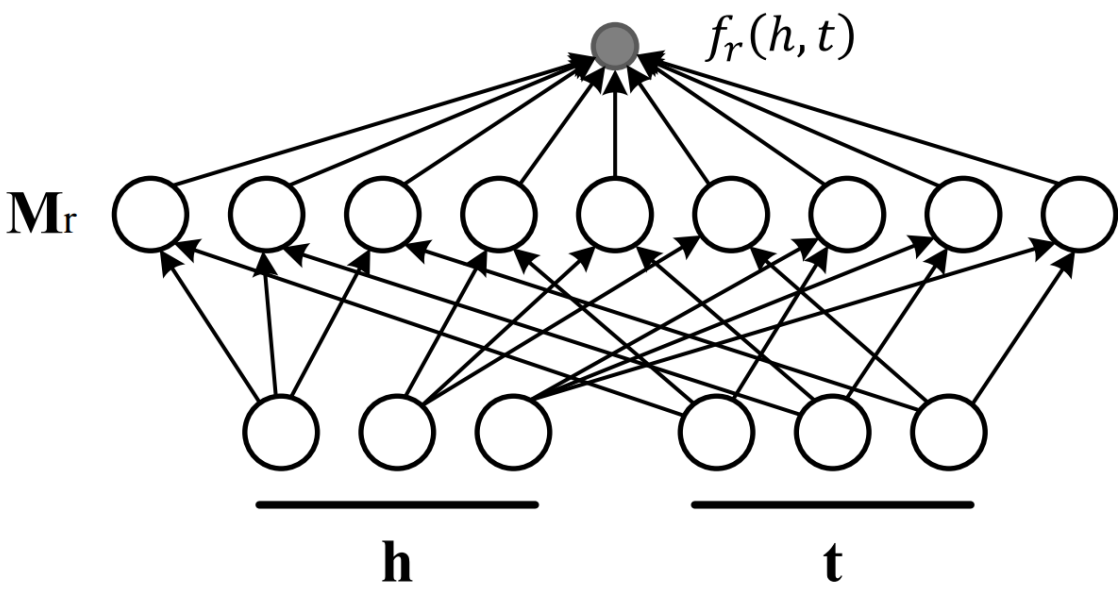
### TransE



Представляет и узлы, и связи как векторы в одном и том же векторном пространстве. Причём оператор отношения рассматривается как вектор перемещения между субъектом и объектом.

$$f_r(h, t) = -\|h + r - t\|_{l/2}$$

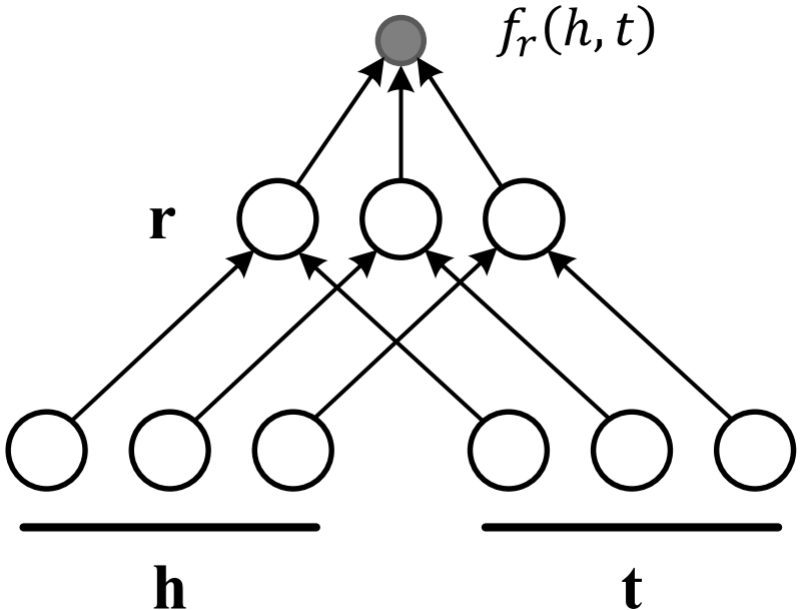
### RESCAL



Ассоциирует каждую сущность с вектором и старается передать скрытые смыслы (факторы) этой сущности. Каждое отношение представлено в виде матрицы, которая моделирует попарные взаимодействия скрытых факторов.

$$f_r(h, t) = h^T M_r t = \sum_{i=0}^{d-1} \sum_{j=0}^{d-1} [M_r]_{ij} \cdot [h]_i \cdot [t]_j$$

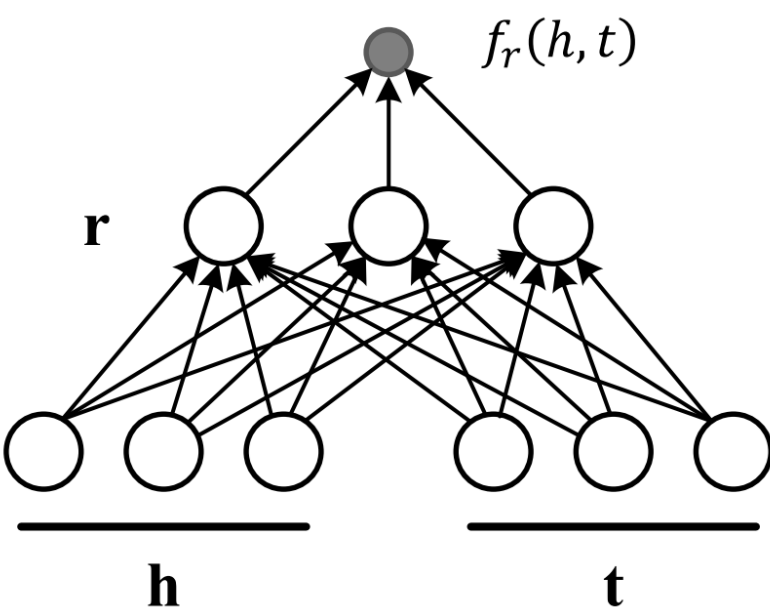
### DistMult



Упрощает модель RESCAL, накладывая ограничение на матрицу отношений: в DistMult Матрица отношений обязательно должна быть диагональной для каждого отношения  $r$

$$f_r(h, t) = h^T \text{diag}(r) t = \sum_{i=0}^{d-1} [r]_i \cdot [h]_i \cdot [t]_i$$

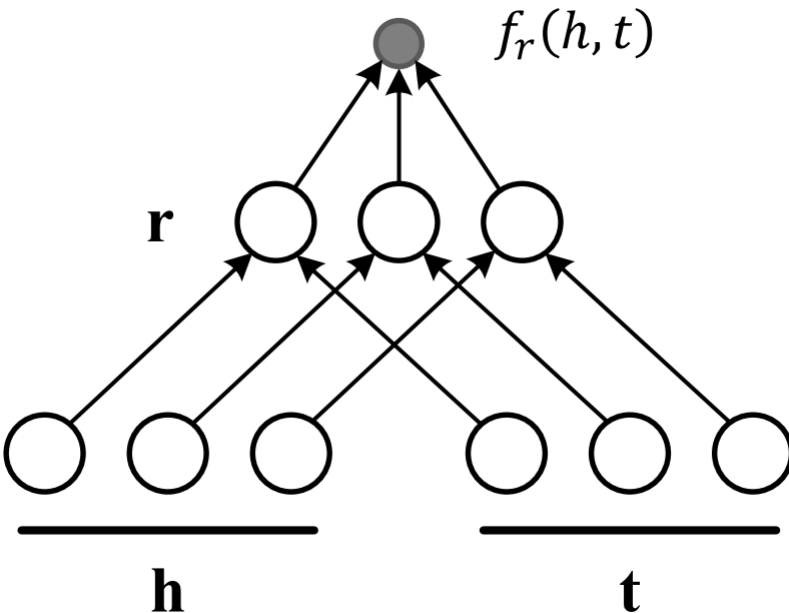
### HolE



Модель голографического эмбеддинга, совмещает в себе мощность RESCAL и простоту DistMult. К представлениям сущностей сначала применяется оператор взаимной корреляции.

$$[h * t]_i = \sum_{k=0}^{d-1} [h]_k \cdot [t]_{(k+i) \bmod d}$$
$$f_r(h, t) = r^T (h * t) = \sum_{i=0}^{d-1} [r]_i \sum_{k=0}^{d-1} [h]_k \cdot [t]_{(k+i) \bmod d}$$

### Complex



Расширяет модель DistMult, производя эмбеддинг в пространство комплексных чисел. За счёт появляется возможность лучше моделировать асимметричные отношения.

$$f_r(h, t) = \text{Re}(h^T \text{diag}(r) \bar{t}) = \text{Re} \left( \sum_{i=0}^{d-1} [r]_i \cdot [h]_i \cdot [\bar{t}]_i \right)$$

Выполнила: Злобина С.В. \_\_\_\_\_

Утвердил: Гапанюк Ю.Е. \_\_\_\_\_

# Подсчёт метрик качества эмбединга

## Задачи для проверки качества эмбединга

### Предсказание отношения

Из истинного триплета  $\langle h, r, t \rangle$  удаляют отношение, получая  $\langle h, ?, t \rangle$ . Задача заключается в предсказании отношения для пары сущностей.

### Предсказание сущности

Из истинного триплета  $\langle h, r, t \rangle$  по очереди удаляют сначала субъект, потом объект, получая 2 неполных триплета:  $\langle ?, r, t \rangle$  и  $\langle h, r, ? \rangle$ . Задача заключается в предсказании недостающих сущностей для пары неполных триплетов.

В модуле KGE выбрана данная задача как наиболее часто встречающаяся

### Классификация триплета

В истинном триплете  $\langle h, r, t \rangle$  заменяют одну из сущностей или отношение так, чтобы он стал ложным. Задача заключается в определении, является ли триплет истинным или ложным.

## Сравнение алгоритмов подсчёта метрик

### Варианты алгоритмов

Обозначение	Алгоритм подсчёта метрик
B1	Mean rank (средний ранг)
B2	Mean reciprocal rank (среднеобратный ранг)
B3	Hit@N

### Критерии оценки алгоритмов

Обозначение	Критерий	Единица измерения
K1	Сложность подсчёта	(качественный критерий)
K2	Удобство в использовании (наглядность)	(качественный критерий)
K3	Показательность на небольших датасетах (качественный критерий)	(качественный критерий)
K4	Устойчивость к выбросам	(качественный критерий)

### Значения критериев

Критерий	Значение	Оценка
Сложность подсчёта	Сложно	1
	Средней сложности	2
	Легко	3
Удобство в использовании (наглядность)	Ненаглядно	1
	Требуется время на привыкание	2
	Наглядно	3
Показательность на небольших датасетах	Совсем непоказателен	1
	Показателен отчасти	2
	Показателен	3
Устойчивость к выбросам	Крайне неустойчив	1
	Средней устойчивости	2
	Устойчив	3

### Оценка вариантов по критериям

все критерии оцениваются по принципу “чем больше, тем лучше”

	Вес критерия	B1	B2	B3
K1	0,1	3	2	2
K2	0,3	1	3	3
K3	0,2	1	2	3
K4	0,4	1	2	3

### Выбор лучшего алгоритма методом взвешенной суммы

	Вес критерия	B1	B2	B3
K1	0,1	1	0,67	0,67
K2	0,3	0,33	1	1
K3	0,2	0,33	0,67	1
K4	0,4	0,33	0,67	1
$\sum_{i=1}^4 (\alpha_i K_i)$		0,397	0,769	0,967

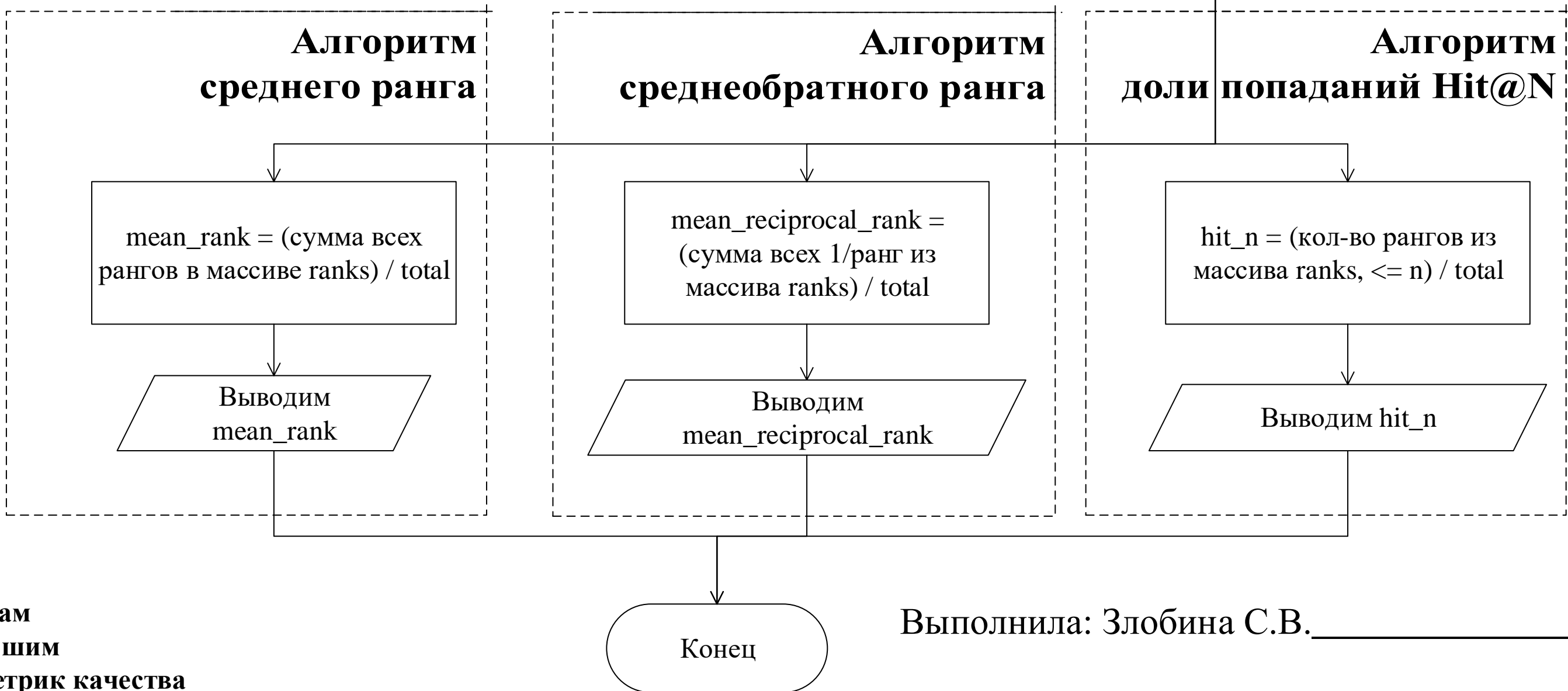
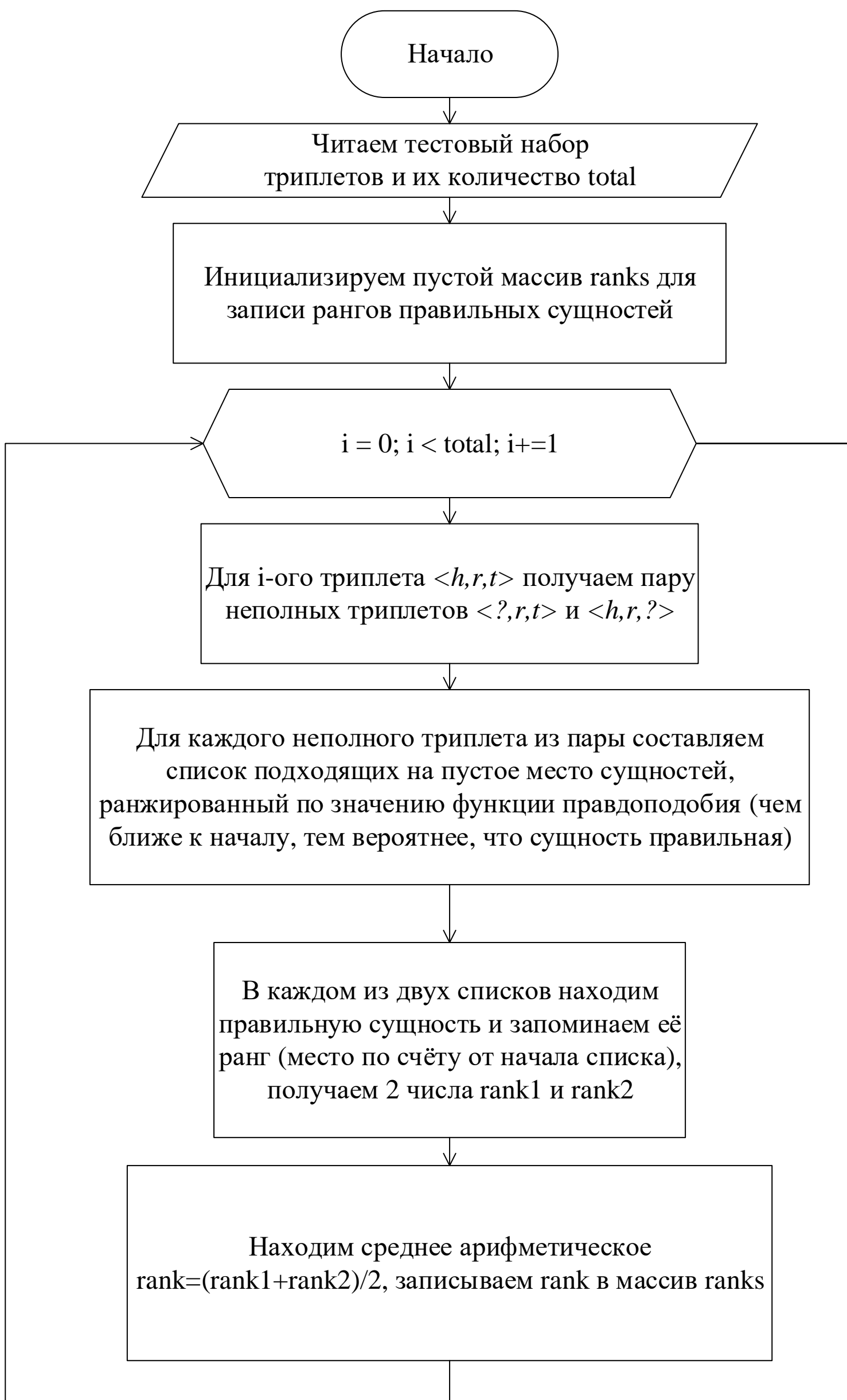
Лучший вариант = Max(взвешенная сумма) = B3

### Выбор лучшего алгоритма методом Борда

	B1	B2	B3
K1	1	2	2
K2	3	1	1
K3	3	2	1
K4	3	2	1
Σ	10	7	5

Лучший вариант = Min(сумма) = B3

## Алгоритмы подсчёта метрик



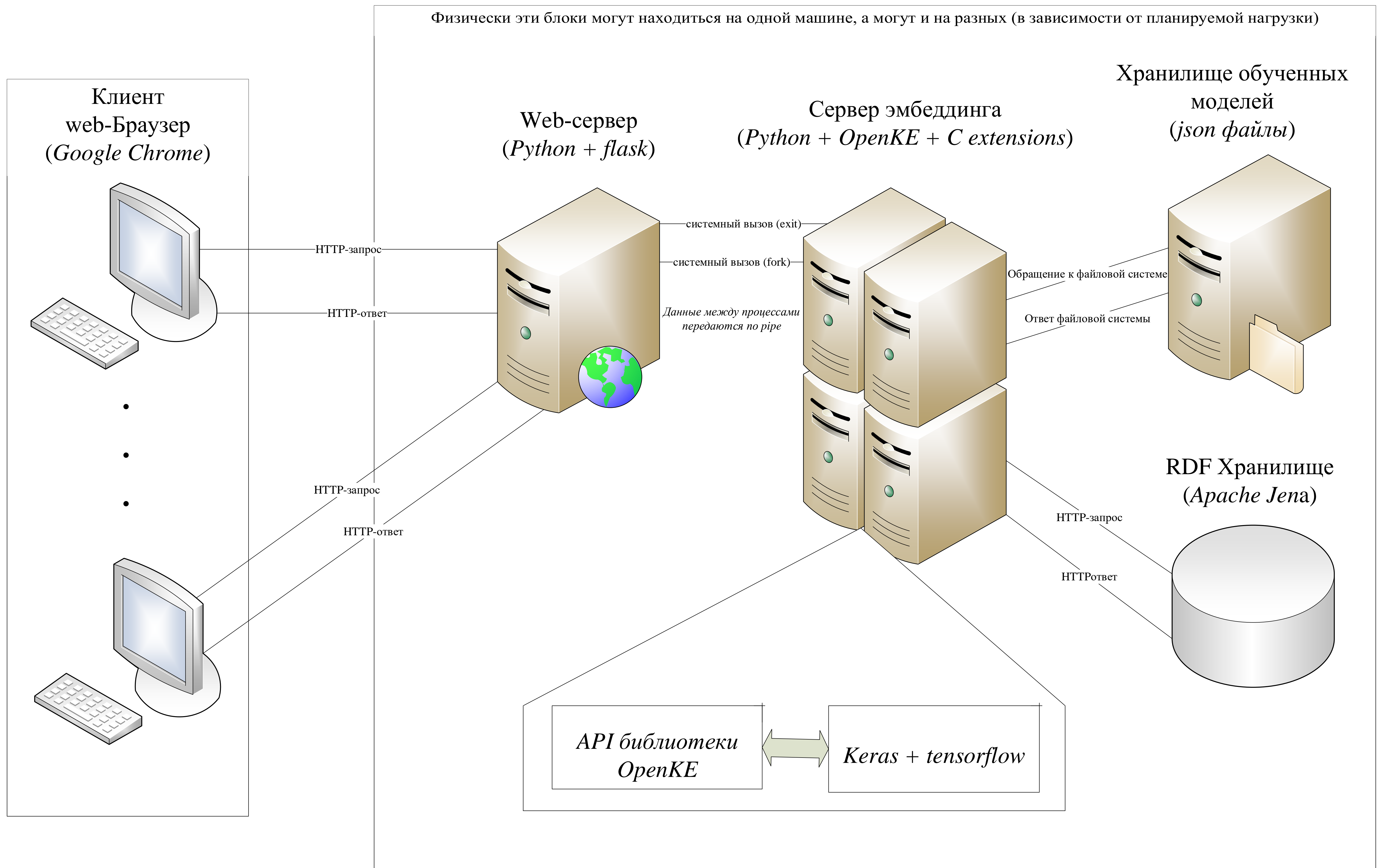
**Вывод:** согласно всем методам проведённых сравнений лучшим алгоритмом для подсчёта метрик качества эмбединга оказался алгоритм B3 (Hit@N), который и был реализован в модуле KGE.

Выполнила: Злобина С.В. \_\_\_\_\_

Утвердил: Гапанюк Ю.Е. \_\_\_\_\_



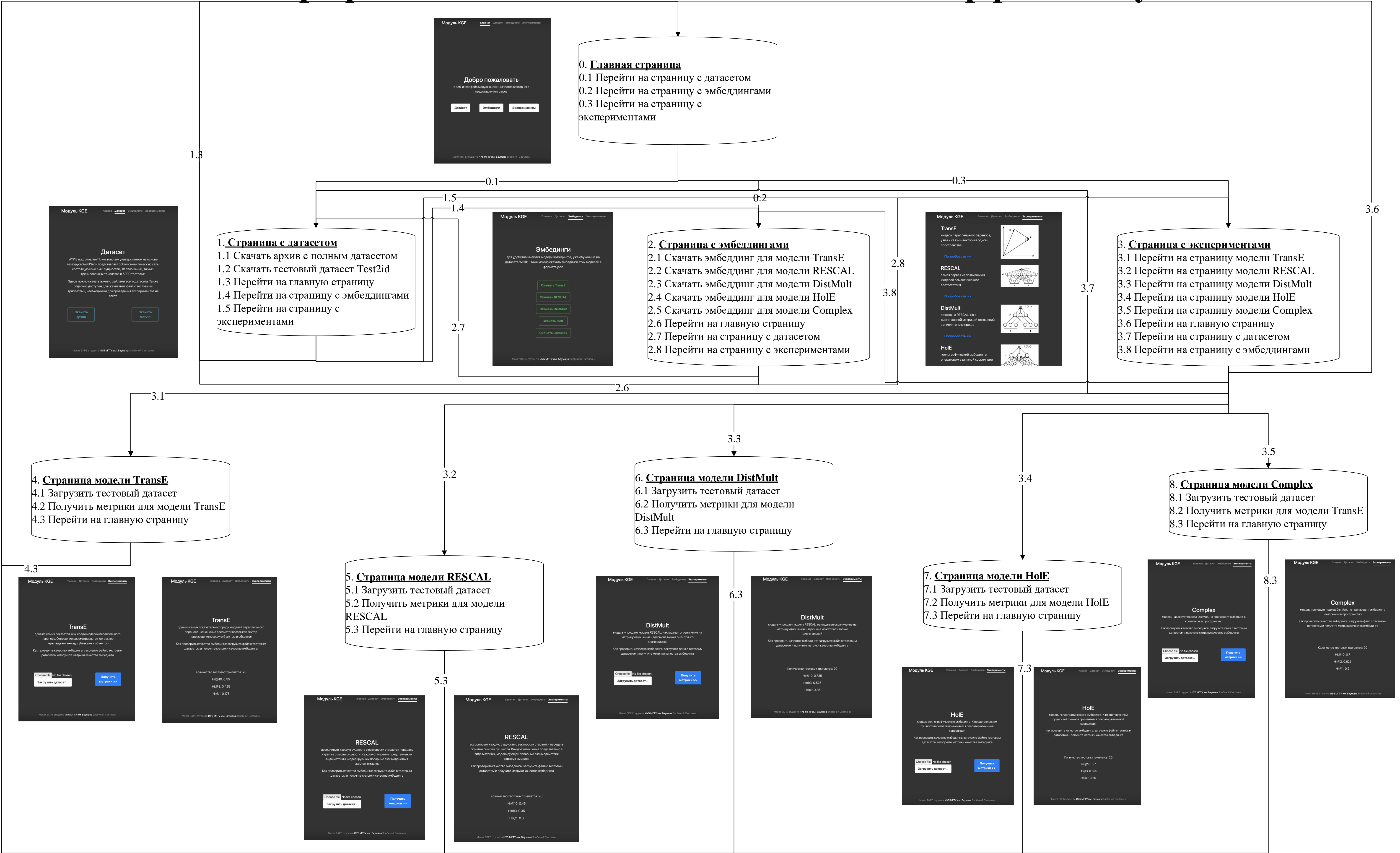
# Архитектура модуля



Выполнила: Злобина С.В. \_\_\_\_\_

Утвердил: Гапанюк Ю.Е. \_\_\_\_\_

# Граф диалога с пользователем и веб-интерфейс модуля



Выполнила: Злобина С.В.\_\_\_\_\_

Утвердил: Гапанюк Ю.Е.\_\_\_\_\_