

# **Обзор основных методов автоматического расширения графов знаний на основе новых данных с точки зрения корректных разрешений возможных конфликтов при слиянии данных.**

Ключевые слова: графы знаний, расширения графов знаний, knowledge graph completion

## **Аннотация**

Данная статья посвящена проблеме расширения графов знаний с точки зрения корректных разрешений возможных конфликтов при слиянии данных на фоне ежегодного увеличения объема генерируемой человечеством информации. В статье выделены ключевые критерии шести основных методов расширения графов знаний для их сравнения, такие как интерпретируемость результатов рассуждений методов расширения, вычислительная сложность и требования к построению этих методов. На основе сравнений методов было выявлено, что существующие решения, обладая высокой степенью интерпретируемости, сталкиваются с проблемами вычислительной сложности, в частности при масштабировании графов знаний, и сложности при построении самих методов. А решения, которые лучше справляются с задачей расширения при масштабировании графов знаний, в большинстве случаев имеют проблемы с интерпретируемостью.

## **Введение**

Объем генерируемых человечеством данных с каждым годом увеличивается: только за последние два года (к 2023 году) было получено примерно 90% мировых данных, а на 2023 год генерация составила 120 зеттабайт [7]. Хранение и структурирование подобной информации в виде графов знаний используются во многих современных информационных технологиях, в частности, в системах автоматической обработки текстов (в том числе, в семантических системах) [8] или в рекомендательных системах, например, в рекомендательной системе для агрегирования междисциплинарных знаний относительно улучшения психического здоровья в графе эмоциональных знаний (EmoKG) [9] или фармакологической рекомендательной системе лекарственных средств [10]. С ростом количества получаемых данных возрастает и сложность автоматического расширения графов знаний с учетом корректного разрешения возможных конфликтов между новыми данными и между новыми и уже имеющимися данными, из-за чего используемые методы обновления и расширения графов знаний могут нуждаться в обновлении. В этой работе предлагается взглянуть на обзор существующих методов расширения графов знаний с точки зрения корректных разрешений возможных конфликтов при слиянии данных. Для этой цели были поставлены следующие задачи:

- Поиск существующих методов автоматического расширения графов знаний
- Определение корректности расширения графов знаний новыми данными для каждого метода с учетом требуемых для него вычислительной сложности и

требований к его построению

## **Обзор предметной области**

### **Принцип отбора аналогов**

Для поиска методов автоматического расширения графов знаний на основе новых данных использовалась система поиска “Академия Google”. Поиск осуществлялся по следующим запросам: <<knowledge graph expansion “knowledge graph expansion”>>, <<knowledge graph completion “knowledge graph completion”>>. “Knowledge graph completion” (KGC) - устоявшийся термин для описания проблемы заполнения графа знаний новыми данными [1]. Все методы заполнения можно разделить на два вида: так называемые, традиционные методы заполнения графов знаний (traditional knowledge graph completion methods) и методы, основанные на глубоком обучении представлений (deep representation learning). ##### Расширение графа знаний на основе логических правил (Rule Reasoning Model) - Традиционный метод Метод заполнения графа знаний на основе логического рассуждения использует правила или статистические характеристики для вывода новых знаний, расширяя структуру графа и дополняя его [2]. Поскольку правила автоматически генерируются в соответствии с семантикой или извлекаются вручную, преимущество метода заключается в его высокой интерпретируемости и точности создаваемых данных в графе знаний. В то же время этот метод также имеет недостатки. Прежде всего, этот метод сильно зависит от правил, которые построить вычислительно трудно, независимо от способа их построения (ручного или автоматического). Причем, при увеличении масштаба графа знаний вычисления новых правил возрастает в разы, из-за чего этот метод становится неприменим. Вычислительная сложность  $O(2 * p * |G|)$ , где  $p$  - количество правил, а  $G$  - количество вершин и связей в графе [11]. ##### Расширение графа знаний на основе вероятностной модели графа (Probabilistic Graph Model) - Традиционный метод Метод заполнения графа знаний на основе вероятностной графовой модели использует графы для представления вероятностных отношений, обеспечивая меньшую вычислительную сложность по сравнению с методами, основанными на правилах [3]. Этот метод преимущественно использует моделирование сетей Маркова и Байесовских сетей (Байесовские сети используются чаще). Сети Маркова используются для представления вероятностных связей, объединяя графовую структуру с теорией вероятностей. Байесовские сети учитывают структуру сети и информацию об атрибутах узлов. Они представляют собой направленный ациклический граф. Преимущества такого подхода: гибкая топологическая структура, высокая интерпретируемость (процесс рассуждений при создании новой связи отслеживается за счет известной структуры представления данных). В то же время она хорошо работает с точки зрения повышения точности прогнозирования и сокращения временных затрат. Однако из-за высокой сложности алгоритма его трудно рассчитать для масштабных графов знаний. Вычислительная сложность  $O(d^3 * r \log d)$ , где  $r$  - глубина в направленном ациклическом графе (максимальное расстояние от корневого узла), а  $d$  - размерность случайного вектора, равного количеству вершин в графе [12]. ##### Расширение графа знаний на основе вычислений графа (Graph Calculation Model) - Традиционный метод Метод заполнения графа знаний на основе графового вычисления

абстрагирует структуру графа знаний в виде графа, где сущности представлены узлами, а отношения различных типов действуют как рёбра [4]. С использованием статистических характеристик узлов и рёбер, таких как степень узла и матрица смежности, можно предсказывать новые сущности и отношения. Этот метод легко поддается интерпретации и не требует дополнительных логических правил, помогающих в процессе рассуждения. Имеет проблемы с масштабируемостью, высокой степенью использования памяти и сталкивается с проблемой сложности крупномасштабных вычислений данных. Вычислительная сложность  $O(|V|^2M)$ , где  $V$  - количество типов связей в графе,  $M$  - количество связей [13]. ##### Расширение графа знаний на основе модели перевода (Translation Model) - Глубокое обучение представлений Расширение графа знаний, основанное на модели трансляции, предсказывает, что новые отношения сущностей утраиваются из существующего графа знаний путем встраивания сущностей и отношений в векторное пространство. Большинство существующих методов фокусируется на структурированной информации троек (субъект, отношение, объект) и максимизируют возможность их установления [5]. Метод заполнения графа знаний, основанный на модели перевода, фокусируется на использовании взаимосвязи между сущностями, семантики, содержащейся в сущности и отношениях, и структурированной информации графа знаний для реализации моделирования сущностей и отношений, что компенсирует сложное обучение и трудное расширение традиционных методов. Для моделирования сущностей и отношений метод имеет высокую интерпретируемость среди всех нетрадиционных методов (в основе метода лежит геометрическая интерпретация векторных пространств - сущности и отношения представляются в виде векторов в пространстве, а операции над векторами используются для представления отношений между сущностями). Вычислительная сложность  $O(n)$  [14] (для TransE, TransH, TransR реализаций). ##### Расширение графа знаний на основе модели семантического сопоставления (Semantic Matching Model) - Глубокое обучение представлений Данный метод использует функцию оценки, основанную на семантической схожести, для извлечения потенциальных семантических ассоциаций между сущностями и отношениями [6]. Путем вложения представлений сущностей и отношений в векторное пространство метод может предсказывать новые факты и расширять граф знаний. Этот метод обеспечивает высокую точность при прогнозировании симметричных отношений, может хорошо выявлять потенциальные семантические ассоциации. В то же время существуют проблемы со многими параметрами модели и высокой вычислительной сложностью, которые не могут быть адаптированы для заполнения крупномасштабных графов знаний. Вычислительная сложность  $O(n^2)$  [14] (для RESCAL реализации). ##### Расширение графа знаний на основе обучения представлений сети (Network Representation Learning Model) - Глубокое обучение представлений Метод, основанный на обучении сетевому представлению, направлен на объединение информации, извлеченной из структуры топологии сети, и информации о содержимом узлов и ребер, преобразование вершин сети во встраиваемые представления в низкоразмерном непрерывном векторном пространстве и реализацию задачи завершения графа знаний с помощью машинного обучения. Применение этого метода, основанного на сетевом представлении, к задаче завершения графа знаний позволяет лучше извлекать скрытые функции в структуре графа знаний, что полезно для обучения модели прогнозирования. Построить нелинейную и разреженную структуру сети

достаточно трудно. Вычислительная сложность  $O(\log n)$  [15] ### Критерии сравнения аналогов

### **Интерпретируемость**

Для автоматического расширения графа знаний на основе новых данных для обработки возможных конфликтов необходимо понимание того, из-за чего они происходят, то есть, интерпретация рассуждений методов расширения, которые привели к соответствующим решениям при добавлении новых сущностей и связей в граф знаний.

### **Сложность вычислений**

Зависимость объёма работы, которая выполняется методом, от размера входных данных.

### **Требования к построению и обучению метода при масштабировании**

Автоматическое расширение графа знаний требует методов, которые могут быть легко настроены при сильно разрастающемся графе знаний.

## **Таблица сравнения по критериям**

Сравнение по критериям представлено в табл. 1.

Таблица 1 – Сравнение

	<b>Интерпретируемость</b>	<b>Сложность вычислений</b>	<b>Требования к построению и обучению метода при масштабировании</b>
<b>Rule Reasoning Model</b>	высокая (так как составляемые правила формулируются на основе логической структуры знаний и отношений в графе - то есть можно проследить, какие шаги и условия привели к конкретному выводу)	$O(2 * p *  G )$	высокие (строить требуемые моделью новые правила, которые бы учитывали все взаимоотношения между вершинами, становится труднее в сильно расрозшемся графе знаний)
<b>Probabilistic Graph Model</b>	высокая ((хоть и более низкая по сравнению с Rule Reasoning Model),	$O(d^3 * r \log d)$	высокие (требует разработки сложных алгоритмов, основанных на моделях

	Интерпретируемость	Сложность вычислений	Требования к построению и обучению метода при масштабировании
	благодаря использованию графов для представления вероятностных отношений (вероятностные выводы основываются на известных графовых вероятностных моделях))		сетей Маркова и Байесовских сетях)
<b>Graph Calculation Model</b>	высокая (поскольку использует графовую структуру для предсказания новых сущностей и отношений)	$O( V ^{2M})$	средние (требует статистических характеристик узлов и рёбер, но не требует дополнительных логических правил)
<b>Translation Model</b>	средняя (в основе метода лежит геометрическая интерпретация векторных пространств - сущности и отношения представляются в виде векторов в пространстве, а операции над векторами используются для представления отношений между сущностями. Интерпретируемость модели, как правило, хуже, чем у традиционных	$O(n)$	средние (требует обширного обучающего набора данных и вычислительных ресурсов для обучения глубоких моделей, однако справляется с крупномасштабными графами знаний)

	Интерпретируемость	Сложность вычислений	Требования к построению и обучению метода при масштабировании
	методов, однако из нетрадиционных методов - это наиболее интерпретируемая модель)		
<b>Semantic Matching Model</b>	низкая (так как вовлекают сложные математические модели)	$O(n^2)$	высокие (не может быть адаптирована для заполнения крупномасштабных графов знаний)
<b>Network Representation Learning Model</b>	низкая (так как вовлекают сложные математические модели)	$O(\log n)$	средние (требует обширного обучающего набора данных и вычислительных ресурсов для обучения глубоких моделей, однако справляется с крупномасштабными графами знаний)

## ***Выводы по итогам сравнения***

Выбор метода для автоматического расширения графа знаний зависит от конкретных требований задачи, доступных ресурсов и важности интерпретируемости результатов. Каждый метод имеет свои преимущества и ограничения, и выбор должен быть обоснован учетом специфики задачи и возможностей инфраструктуры. В контексте задачи расширения графов знаний новыми данными, где возможны конфликты при слиянии и необходимости их обрабатывать, наиболее важным критерием становится интерпретируемость выводов выбранного метода, что в большей мере свойственно традиционным методам. Тем не менее, существующие традиционные методы сталкиваются с проблемой масштабируемости графов знаний, что на фоне ежегодного увеличения объема генерируемых человечеством данных [7] становится существенной проблемой.

## ***Выбор метода решения***

Метод решения задачи расширения графов знаний новыми данными, где возможны конфликты при слиянии и необходимости их обрабатывать, должен достигать такой же

высокой интерпретируемости, так как и традиционные методы, так как важно понимание того, какие знания были добавлены и какие выводы сделаны. Но при этом быть более пригодным для крупномасштабных графов, потому как с каждым годом увеличивается объем генерируемых человечеством данных, с чем традиционные модели справляются плохо.

## **Описание метода решения**

На данный момент существует два наиболее подходящих варианта решения задачи расширения графов знаний новыми данными, где возможны конфликты при слиянии и необходимости их обрабатывать: либо использовать один из традиционных методов, потому как они обладают высокой степенью интерпретируемости, но низкой степенью масштабируемости, либо же использовать метод Translation Model, который хоть и не обладает той же высокой степенью интерпретируемости, однако лучше работает на крупномасштабных графах.

## **Заключение**

Обзор основных методов расширения графов на основе новых данных с точки зрения корректных разрешений возможных конфликтов при слиянии данных показал, что существующие решения могут справляться с данной задачей на небольших данных, но при масштабировании графа знаний все традиционные методы сталкиваются с проблемой вычислительной сложности из-за специфики построения метода. Методы, которые лучше работают при масштабировании, напротив же, зачастую имеют проблемы с интерпретируемостью, так как вовлекают более сложные математические модели.

В качестве дальнейших исследований следует более подробно сравнить Translation Model с традиционными методами, например для конкретных реализаций. Также в этой статье затрагивались только основные методы, помимо них существуют и другие методы, такие как ConvE, которые лучше могут справляться с рассматриваемой задачей. ## Список литературы

- [1] Chen Z. et al. Knowledge graph completion: A review //Ieee Access. – 2020. – Т. 8. – С. 192435-192456.
- [2] Carlson A. et al. Toward an architecture for never-ending language learning //Proceedings of the AAAI conference on artificial intelligence. – 2010. – Т. 24. – №. 1. – С. 1306-1313.
- [3] Jiang S., Lowd D., Dou D. Learning to refine an automatically extracted knowledge base using markov logic //2012 IEEE 12th International Conference on Data Mining. – IEEE, 2012. – С. 912-917.
- [4] Lao N., Cohen W. W. Relational retrieval using a combination of path-constrained random walks //Machine learning. – 2010. – Т. 81. – С. 53-67.
- [5] Wang Q. et al. Knowledge graph embedding: A survey of approaches and applications //IEEE

- Transactions on Knowledge and Data Engineering. – 2017. – T. 29. – №. 12. – C. 2724-2743.
- [6] Yang B. et al. Embedding entities and relations for learning and inference in knowledge bases //arXiv preprint arXiv:1412.6575. – 2014.
- [7] Hassani H., MacFeely S. Driving Excellence in Official Statistics: Unleashing the Potential of Comprehensive Digital Data Governance //Big Data and Cognitive Computing. – 2023. – T. 7. – №. 3. – C. 134.
- [8] Schneider P. et al. A decade of knowledge graphs in natural language processing: A survey //arXiv preprint arXiv:2210.00105. – 2022.
- [9] Gyrard A., Boudaoud K. Interdisciplinary iot and emotion knowledge graph-based recommendation system to boost mental health //Applied Sciences. – 2022. – T. 12. – №. 19. – C. 9712.
- [10] Sosa D. N. et al. A literature-based knowledge graph embedding method for identifying drug repurposing opportunities in rare diseases //Pacific Symposium on Biocomputing 2020. – 2019. – C. 463-474.
- [11] Ahmadi N. et al. Mining expressive rules in knowledge graphs //Journal of Data and Information Quality (JDIQ). – 2020. – T. 12. – №. 2. – C. 1-27.
- [12] Gao M., Aragam B. Efficient Bayesian network structure learning via local Markov boundary search //Advances in Neural Information Processing Systems. – 2021. – T. 34. – C. 4301-4313.
- [13] Lao N. Efficient random walk inference with knowledge bases : дис. – Carnegie Mellon University, 2012.
- [14] Ebisu T., Ichise R. Toruse: Knowledge graph embedding on a lie group //Proceedings of the AAAI conference on artificial intelligence. – 2018. – T. 32. – №. 1.
- [15] Perozzi B., Al-Rfou R., Skiena S. Deepwalk: Online learning of social representations //Proceedings of the 20th ACM SIGKDD international conference on Knowledge discovery and data mining. – 2014. – C. 701-710.
- [16] Dettmers T. et al. Convolutional 2d knowledge graph embeddings //Proceedings of the AAAI conference on artificial intelligence. – 2018. – T. 32. – №. 1.