

# 知识图谱推理问答综述

清华大学知识工程实验室 史佳欣

# | 讲者介绍



- 清华大学知识工程实验室
- 博士四年级
- 导师为李涓子教授
- 研究兴趣包括视觉推理，知识图谱推理，文档摘要等

# | 什么是推理问答



- 简单问答涉及单个实体和单个关系

中国的首都是哪里？

What is the capital of China?

- **推理问答**的问题相对复杂，常常涉及多个实体，多个关系，多跳，比较等

中国的首都和美国的首都，哪个人口更多？

Does the capital of China or the capital of America have more population?

# I 推理问答的难点



- 简单问答只需要识别出问题中的实体和关系，链接到知识图谱中，即可查出答案
- 推理问答要求计算机具备多种推理能力，具体包括：
  - 处理多跳关系的能力，如“姚明的妻子的学校”
  - 数值比较的能力，如“哪个城市的人口更多”
  - 集合操作的能力，如“即是篮球运动员，又是球队老板的人有哪些”
  - .....



# 数据集

# 知识图谱推理问答的数据集



数据集	知识库	知识类型	问题数量	自然语言	SPARQL
LC-QuAD2.0 [1]	Wikidata and DBpedia	多种	30k	是	有
ComplexWebQuestions [2]	Freebase	多种	35k	是	有
MetaQA [3]	WikiMovies	单种	400k	否	无
CSQA [4]	Wikidata	单种	1.6M	否	无

# | 知识图谱推理问答的数据集



- 现有知识库的三种知识类型：
  - 关系型，如（“姚明”，“出生于”，“上海”）
  - 属性型，如（“姚明”，“身高”，“229 厘米”）
  - 事实型，用于表示一个关系型事实或属性型事实的知识，如（（“上海”，“人口”，“23,390,000”），“统计时间”，“2016”）
- MetaQA 和 CSQA 仅考虑关系型知识

# | 知识图谱推理问答的数据集



- 现有数据集都缺乏推理过程
- 人类是如何学习解答复杂问题的
  - 先学会解答简单问题
  - 再学会将复杂问题分解为简单问题的组合
- 如果数学老师只讲答案会怎么样?



# 知识图谱推理问答的数据集



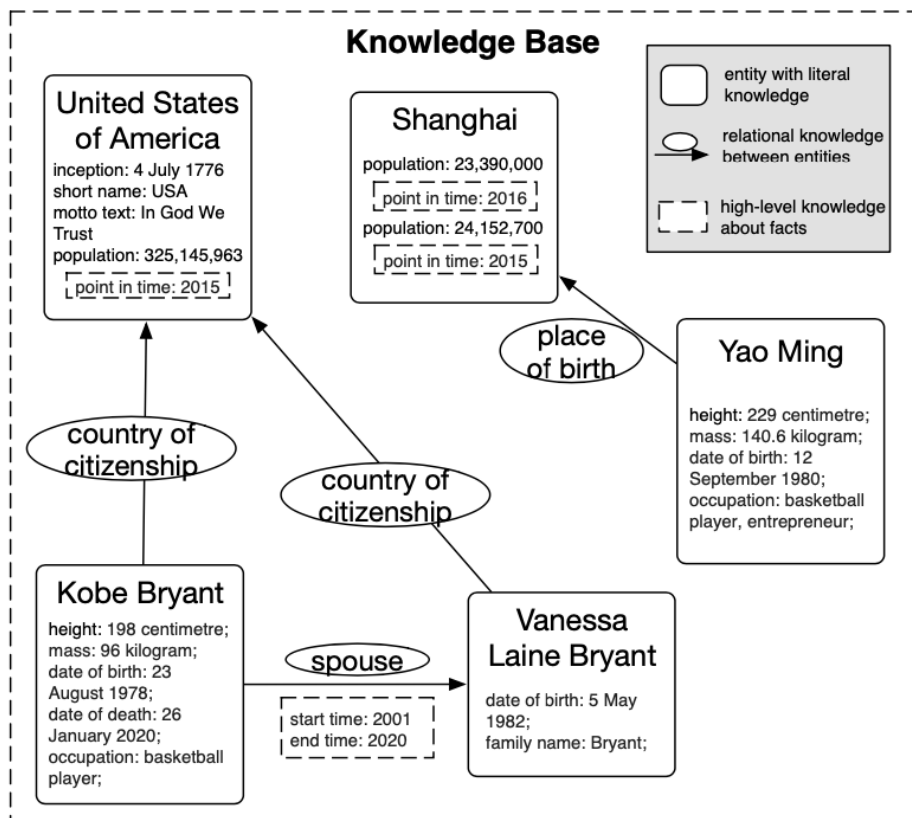
数据集	知识库	知识类型	问题数量	自然语言	SPARQL	推理过程
LC-QuAD2.0 [1]	Wikidata and DBpedia	多种	30k	是	有	无
ComplexWebQuestions [2]	Freebase	多种	35k	是	有	无
MetaQA [3]	WikiMovies	单种	400k	否	无	无
CSQA [4]	Wikidata	单种	1.6M	否	无	无
<b>KQA Pro (Ours)</b>	<b>Wikidata子集</b>	<b>多种</b>	<b>120k</b>	<b>是</b>	<b>有</b>	<b>有</b>

# | 知识图谱推理问答的数据集



- 如何表示推理过程
- 函数 (function) 对应简单问题
- 程序 (program) 对应复杂问题
- 程序由函数组合而成

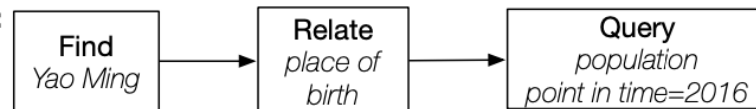
# 知识图谱推理问答的数据集



**Question 1:** How many people does the Yao Ming's birth place have in 2016?

**SPARQL:** `SELECT ?v WHERE { ?e_1 <name> "Yao Ming" . ?e_1 <place_of_birth> ?e . ?e <population> ?v . [ <fact_h> ?e ; <fact_r> <population> ; <fact_t> ?v ] <point_in_time> 2016 . }`

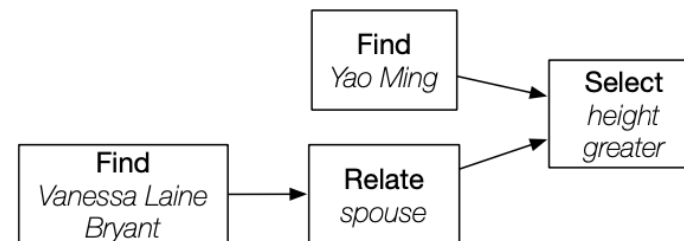
**Program:**



**Question 2:** Who is taller, Yao Ming or the spouse of Vanessa Laine Bryant?

**SPARQL:** `SELECT ?e WHERE { { ?e <name> "Yao Ming" . } UNION { ?e_1 <name> "Vanessa Laine Bryant" . ?e_1 <spouse> ?e . } ?e <height> ?v . } ORDER BY DESC(?v) LIMIT 1`

**Program:**





# 方法

# | 知识图谱推理问答的方法



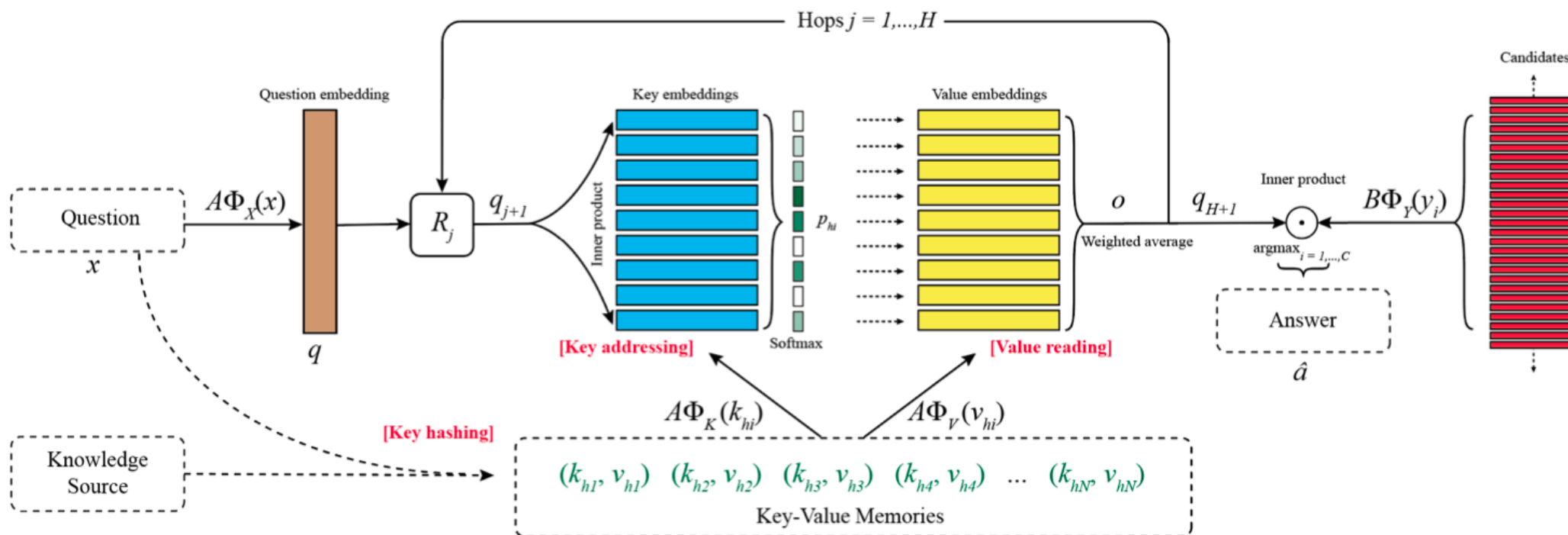
- 键值记忆网络 (KVMemNet)
- 基于强化学习的多跳路径搜索
- 弱监督的程序学习
- 查询图解析与匹配

# | 知识图谱推理问答的方法



- 键值记忆网络 (KVMemNet)
- 基于强化学习的多跳路径搜索
- 弱监督的程序学习
- 查询图解析与匹配

# 键值记忆网络



KVMemNet 框架图

# I 键值记忆网络



- Key Hashing: 将知识库转换成  $(k, v)$  的形式, 并从中选取一个子集
  - 形式转换: 对于三元组  $(s, p, o)$ , 将  $s$  和  $p$  共同作为  $k$ , 将  $o$  作为  $v$
  - 子集选择条件
    - $k$  与输入的问题有共同的单词
    - 共同的单词不是停用词
    - 根据共同的单词数量排序, 选择前  $N$  个
- 每个问题都需要构造对应的 Memory



# I 键值记忆网络



- Key Addressing: 根据问题  $x$ , 为 Memory 中的所有 key 计算一个概率分布
- Value Reading:  $p_{h_i} = \text{Softmax}(A\Phi_X(x) \cdot A\Phi_K(k_{h_i}))$ :

$$o = \sum_i p_{h_i} A\Phi_V(v_{h_i})$$

# | 键值记忆网络



- Query Updating: 用得到的 value 向量更新 query 向量, 使用映射矩阵  $R_j$
- 用  $q_{j+1}$  替换 Key Addressing 中的问题向量, 迭代更新
- 迭代  $H$  步之后, 将  $q_{H+1}$  输入分类器中, 预测答案

$$p_{h_i} = \text{Softmax}(q_{j+1}^\top A \Phi_K(k_{h_i}))$$

# I 键值记忆网络



- 优点：
  - 模型简单，通用性强
  - 通过向量的迭代更新，隐式进行推理
- 缺点：
  - 需要对每个问题构造 Memory，容易占用大量的时间和空间
  - 推理能力较弱
  - 缺乏可解释性

# | 知识图谱推理问答的方法



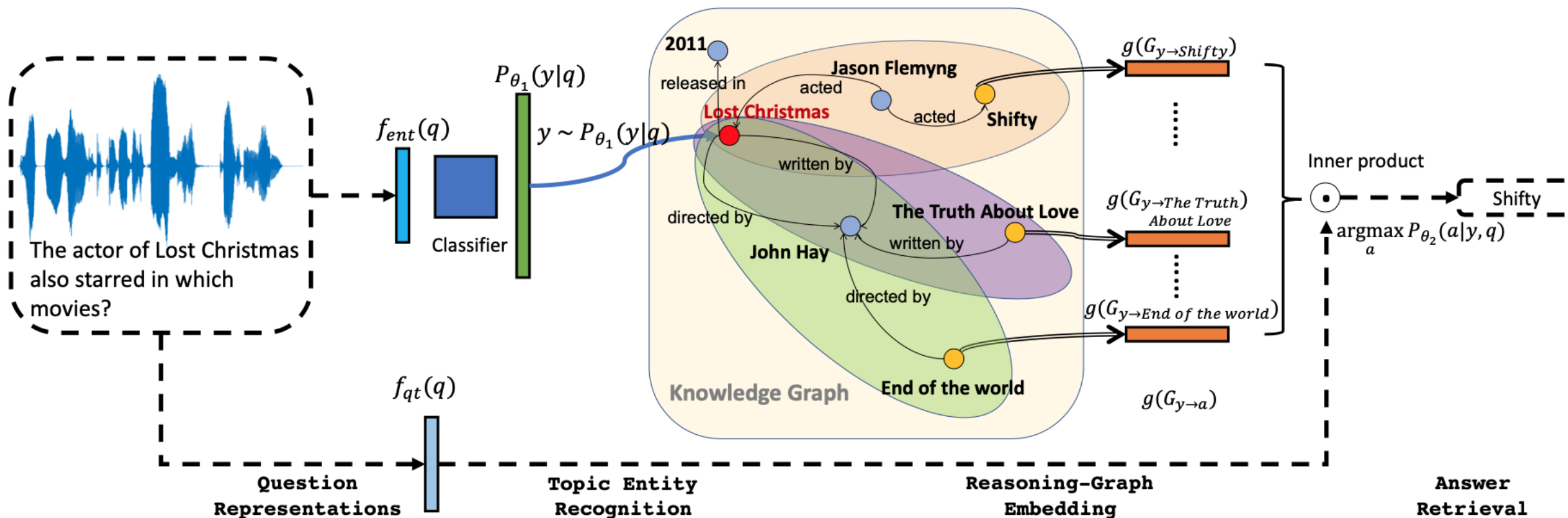
- 键值记忆网络 (KVMemNet)
- 基于强化学习的多跳路径搜索
- 弱监督的程序学习
- 查询图解析与匹配

# | 基于强化学习的多跳路径搜索



- 基本思路：
  - 找到问题中的主题实体（Topic Entity），并链接到知识库上
  - 从知识库上的主题实体出发，根据问题选择一个关系，从而跳转到一个新的实体
  - 继续选择关系，跳转实体
  - 迭代若干步，用最终的实体作为答案
  - 利用强化学习进行训练

# Variational Reasoning Network



Variational Reasoning Network (VRN) 框架图

# | Variational Reasoning Network



分为两个模块:

- 主题实体识别模块  $P_{\theta_1}(y|q)$
- 知识推理模块  $P_{\theta_2}(a|y, q)$ , 将主题实体  $y$  看作隐变量

两个模块进行联合优化:

$$\max_{\theta_1, \theta_2} \frac{1}{N} \sum_{i=1}^N \log \left( \sum_{y \in V(\mathcal{G})} P_{\theta_1}(y|q_i) P_{\theta_2}(a_i|y, q_i) \right).$$

# | Variational Reasoning Network



主题实体识别模块:

$$\begin{aligned} P_{\theta_1}(y|q) &= \text{softmax} \left( W_y^\top f_{\text{ent}}(q) \right) \\ &= \frac{\exp(W_y^\top f_{\text{ent}}(q))}{\sum_{y' \in V(\mathcal{G})} \exp(W_{y'}^\top f_{\text{ent}}(q))}, \end{aligned}$$

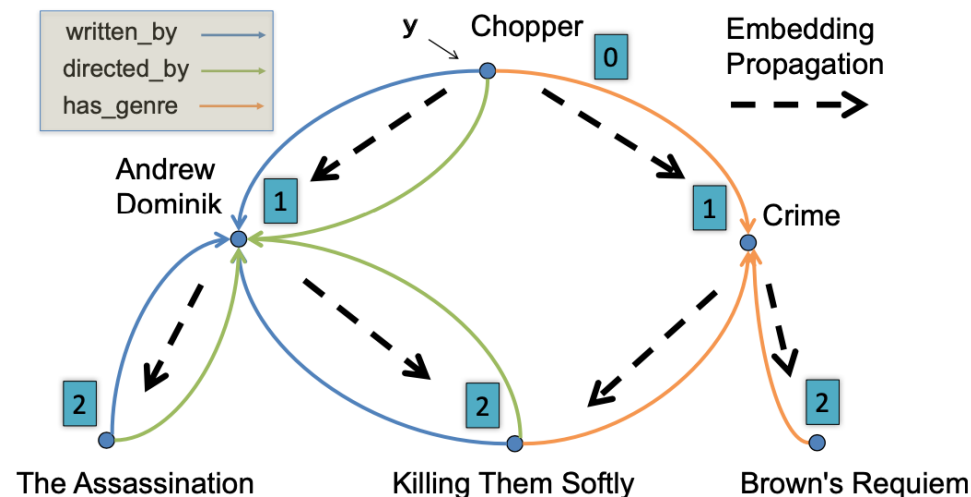


# Variational Reasoning Network



知识推理模块:

- 从主题实体出发, 找出所有能在  $T$  跳之内到达的实体 (忽略边的方向) 及对应的整个子图, 用  $G_y$  表示



# Variational Reasoning Network



知识推理模块:

- $G_y$  中包含的实体都是潜在的答案实体
- 对任一潜在答案  $a$ , 定义  $G_{y \rightarrow a}$  为包含了所有从  $y$  到  $a$  路径的子图
- 计算子图的向量表示  $g(G_{y \rightarrow a}) = \frac{1}{\#\text{Parent}(a)} \sum_{a_j \in \text{Parent}(a), (a_j, r, a) \text{ or } (a, r, a_j) \in G_y} \sigma(V \times [g(G_{y \rightarrow a_j}), \vec{e}_r]),$
- 将  $g(G_{y \rightarrow a})$  作为  $a$  的特征向量用于计算概率

$$\begin{aligned} P_{\theta_2}(a|y, q) &= \text{softmax} \left( f_{\text{qt}}(q)^\top g(G_{y \rightarrow a}) \right) \\ &= \frac{\exp(f_{\text{qt}}(q)^\top g(G_{y \rightarrow a}))}{\sum_{a' \in V(G_y)} \exp(f_{\text{qt}}(q)^\top g(G_{y \rightarrow a'}))}. \end{aligned}$$

# Variational Reasoning Network



训练:

- 目标函数的原始形式为  $\max_{\theta_1, \theta_2} \frac{1}{N} \sum_{i=1}^N \log \left( \sum_{y \in V(\mathcal{G})} P_{\theta_1}(y|q_i) P_{\theta_2}(a_i|y, q_i) \right).$
- 使用辨分推断得到目标函数的下界

$$\begin{aligned} \max_{\psi, \theta_1, \theta_2} \mathcal{L}(\psi, \theta_1, \theta_2) = & \frac{1}{N} \sum_{i=1}^N \mathbb{E}_{Q_{\psi}(y|q_i, a_i)} [ \\ & \log P_{\theta_1}(y|q_i) + \log P_{\theta_2}(a_i|y, q_i) \\ & - \log Q_{\psi}(y|q_i, a_i)], \end{aligned}$$

- 用强化学习进行优化

# I 基于强化学习的多跳路径搜索



- 优点：
  - 有较强的多跳推理能力
  - 可以得到推理路径，具有较好的可解释性
- 缺点：
  - 只能处理关系型知识，无法处理属性型、事实型
  - 问题中必须有且仅有一个主题实体，适用范围较小
  - 不具备其他推理能力

# | 知识图谱推理问答的方法



- 键值记忆网络 (KVMemNet)
- 基于强化学习的多跳路径搜索
- 弱监督的程序学习
- 查询图解析与匹配

# | 弱监督的程序学习



- 基本思路：
  - 定义一些基本的函数，每个函数负责特定功能，用规则实现
  - 将输入问题解析为程序（即函数的组合），执行程序得到答案
  - 由于缺乏监督信息，已有方法使用如下方式学习程序：
    - 强化学习 [1,3,4]，跑出正确答案的程序获得正向激励 **难点：搜索空间巨大，难以收敛**
    - 枚举程序并执行，将跑出正确答案的程序作为标注 [2]

[1] Chen Liang, Jonathan Berant, Quoc Le, Kenneth Forbus, and Ni Lao. Neural symbolic machines: Learning semantic parsers on freebase with weak supervision. In ACL, 2017.

[2] Daya Guo, Duyu Tang, Nan Duan, Ming Zhou, and Jian Yin. Dialog-to-action: Conversational question answering over a large-scale knowledge base. In Advances in Neural Information Processing Systems, 2018.

[3] Amrita Saha, Ghulam Ahmed Ansari, Abhishek Laddha, Karthik Sankaranarayanan, and Soumen Chakrabarti. Complex program induction for querying knowledge bases in the absence of gold programs. Transactions of the Association for Computational Linguistics, 2019.

[4] Ghulam Ahmed Ansari, Amrita Saha, Vishwajeet Kumar, Mohan Bhambhani, Karthik Sankaranarayanan, and Soumen Chakrabarti. Neural program induction for kbqa without gold programs or query annotations. In IJCAI, 2019.

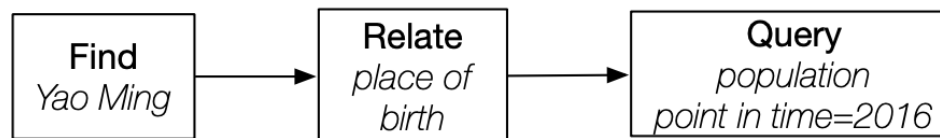
# | 弱监督的程序学习



**Question 1:** How many people does the Yao Ming's birth place have in 2016?

**SPARQL:** `SELECT ?v WHERE { ?e_1 <name> "Yao Ming" . ?e_1 <place_of_birth> ?e . ?e <population> ?v . [ <fact_h> ?e ; <fact_r> <population> ; <fact_t> ?v ] <point_in_time> 2016 . }`

**Program:**

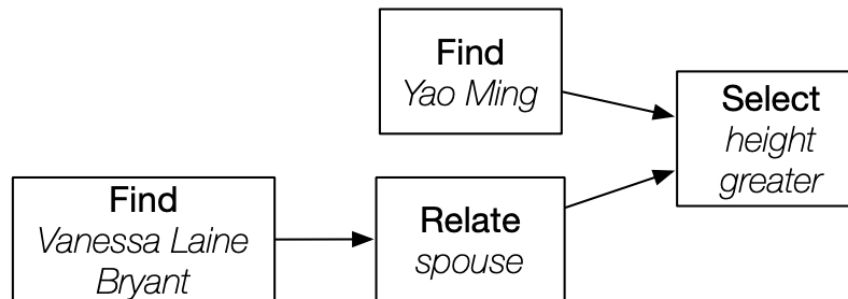


## KQA Pro 程序示例

**Question 2:** Who is taller, Yao Ming or the spouse of Vanessa Laine Bryant?

**SPARQL:** `SELECT ?e WHERE { { ?e <name> "Yao Ming" . } UNION { ?e_1 <name> "Vanessa Laine Bryant" . ?e <spouse> ?e_1 . } ?e <height> ?v . } ORDER BY DESC(?v) LIMIT 1`

**Program:**



# | 弱监督的程序学习



- 引入约束和 trick 以帮助收敛：
  - 函数的语法约束，如函数 A 必须在 B 之后
  - 函数的位置约束，如程序必须以 A 或 B 开头，以 C 或 D 结尾
  - 辅助激励策略，如程序给出的答案类型正确时，给予一定的正向激励
  - 根据目标答案的类型进行剪枝
  - .....



# I 弱监督的程序学习



- 优点：
  - 只需定义相应的函数，理论上可以处理任何推理问题
  - 组合性，有限函数的组合可以解决无限的问题
  - 程序表示推理步骤，具有很好的可解释性
- 缺点：
  - 搜索空间巨大，优化困难，非常耗时，性能不理想
  - 需要人工定义函数的具体实现，不易扩展

# | 知识图谱推理问答的方法



- 键值记忆网络 (KVMemNet)
- 基于强化学习的多跳路径搜索
- 弱监督的程序学习
- 查询图解析与匹配

# | 查询图解析与匹配



- 基本思路:
  - 将自然语言问题解析为查询图的形式，如使用翻译模型将问题转换为 SPARQL
  - 用查询图与知识图谱进行匹配，找出答案

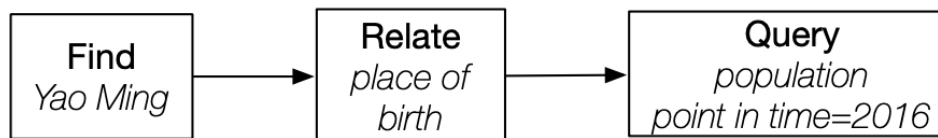
# | 查询图解析与匹配



**Question 1:** How many people does the Yao Ming's birth place have in 2016?

**SPARQL:** `SELECT ?v WHERE { ?e_1 <name> "Yao Ming" . ?e_1 <place_of_birth> ?e . ?e <population> ?v . [ <fact_h> ?e ; <fact_r> <population> ; <fact_t> ?v ] <point_in_time> 2016 . }`

**Program:**

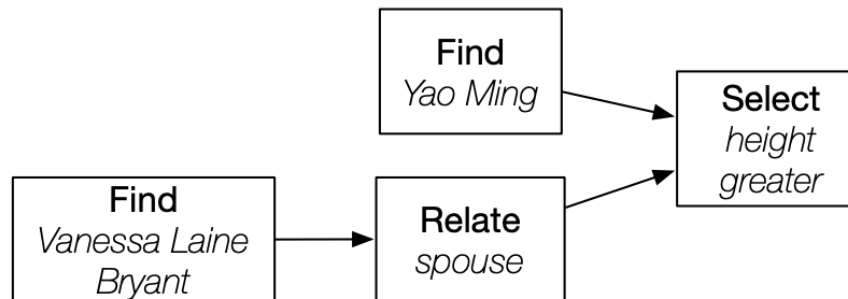


## KQA Pro SPARQL 示例

**Question 2:** Who is taller, Yao Ming or the spouse of Vanessa Laine Bryant?

**SPARQL:** `SELECT ?e WHERE { { ?e <name> "Yao Ming" . } UNION { ?e_1 <name> "Vanessa Laine Bryant" . ?e <spouse> ?e_1 . } ?e <height> ?v . } ORDER BY DESC(?v) LIMIT 1`

**Program:**



# | 查询图解析与匹配



- 优点：
  - 可以处理大多数推理问题
  - 查询图具有较好的可解释性
- 缺点：
  - 由于查询图的复杂性，解析的准确率往往比较低
  - 需要大量的训练数据



# 展望

# | 未来研究方向



- 基于 KQA Pro 数据集
- 神经模块网络, neural-symbolic
  - 将函数实现为神经网络, 以获得更好的鲁棒性
- sequence-to-graph 模型进行查询图解析
  - 将 SPARQL 看作图而非序列, 以更好地捕获节点间的依赖关系
- 用于推理任务的 GCN、RGCN 模型

# Q & A