FAME: Flexible, Scalable Analogy Mappings Engine

Anonymous ACL submission

Abstract

001Analogy is one of the core capacities of human
cognition; when faced with new situations, we
often transfer prior experience from other do-
mains. Most work on computational analogy
relies heavily on complex, manually crafted in-
put. In this work, we relax the input require-
ments, requiring only names of entities to be
mapped. We automatically extract common-
sense representations and use them to identify
a mapping between the entities. Unlike previ-
ous works, our framework can handle partial
analogies, suggesting new entities to be added.
Our method's output is easily interpretable.

Experiments show that our model correctly maps 81.2% of classical 2x2 analogy problems. On larger problems, it achieves 77.8% accuracy (mean guess level: 13.1%). In another experiment, we show our algorithm outperforms human performance, and the automatic suggestions of new entities resemble those suggested by humans. We hope this work will advance computational analogy by paving the way to more flexible, realistic input requirements, with broader applicability.

1 Introduction

016

017

022

024

037

One of the pinnacles of human cognition is the ability to find parallels across distant domains and transfer ideas between them. This *analogous reasoning* process enables us to learn new information faster and solve problems based on prior experience (Minsky, 1988; Hofstadter and Sander, 2013; Holyoak, 1984; PJM, 1966).

The most seminal work in computational analogy is Gentner's Structure Mapping Theory (SMT) (Gentner, 1983) and its implementation, Structure Mapping Engine (SME) (Falkenhainer et al., 1989). In a nutshell, SMT maps between objects in a base domain and objects in a target domain. The mapping is based on a common *relational structure*, rather than on object attributes. For example, consider the Rutherford model of the hydrogen atom, where the atom was explained in terms of the (better-understood) solar system (Falkenhainer et al., 1989). A planet revolving around the sun is mapped to an electron revolving around the nucleus. The mapping is due to shared *relations* between objects (revolving around, being attracted to), not object attributes (round, small).

041

042

043

044

045

047

049

052

054

058

059

060

061

062

063

064

065

066

067

068

069

070

071

074

075

076

077

078

079

One of the main criticisms brought against SME and its follow-up work is their need for extensive hand-coded input – structured representations of both the entities and their relations (see Figure 1 for the input to the atom/solar system mapping). Chalmers et al. (1992) argued that too much human creativity is required to construct this input, and the analogy is already effectively given in the representations: "A brief examination [...] shows that the discovery of the similar structure in these representations is not a difficult task. The representations have been set up in such a way that the common structure is immediately apparent. Even for a computer program, the extraction of such common structure is relatively straightforward."

Some follow-up works avoid hand-coding LISPlike representations, generating them from sketches (Forbus et al., 2011), qualitative simulators (Dehghani and Forbus, 2009), etc. However, they still require much knowledge engineering, and thus are hard to scale. Nowadays, when the web is full of information about potential domains to transfer ideas from (McNeil Jr and Odón, 2013), such representations do not tap into the potential of web-scale analogies for augmenting human creativity.

The method with the simplest input we are aware of is Latent Relation Mapping Engine (LRME) (Turney, 2008), which requires only two lists of entities to be mapped. Given two entities, they search a large corpus for phrases containing both and use them to generate patterns. For example, "a sun centered solar system illustrates" gives rise to patterns such as "a X * Y illustrates". However, such pat-

Solar System	
Solar System (defEntity sun :type inanimate) (defEntity planet :type inanimate) (defDescription solar-system entities (sun planet) expressions (((mass sun) :name mass-sun) ((mass planet) :name mass-planet) ((greater mass-sun mass-planet) :name >mass) ((attracts sun planet) :name attracts) ((revolve-around planet sun) :name revolve) ((and >mass attracts) :name andt) ((cause andt revolve) :name cause-revolve) ((temperature sun) :name temp-planet) ((temperature planet) :name temp-planet) ((greater temp-sun temp-planet) :name >temp)	Rutherford atom (defEntity nucleus :type inanimate) (defEntity electron :type inanimate) (defDescription rutherford-atom entities (nucleus electron) expressions (((mass nucleus) :name mass-n) ((mass electron) :name mass-e) ((greater mass-n mass-e) : name >mass) ((attracts nucleus electron) :name attracts) ((revolve-around electron nucleus) :name revolve) ((charge electron) :name q-electron) ((charge nucleus) :name q-nucleus) ((opposite-sign q-nucleus q-electron) :name >charge)
((greater temp-sun temp-planet) :name >temp) ((gravity mass-sun mass-planet) :name force-gravity) ((cause force-gravity attracts) :name why-attracts)))	((cause >charge attracts) :name why-attracts)))

Figure 1: SME representation of the Solar system/Rutherform atom. Reproduced from Falkenhainer et al. (1989).

terns are extremely simple and brittle, and LRME requires exact matches between the domains (so "revolve around" is different from "rotate around").

In this work, we develop FAME, a Flexible Analogy Mapping Engine. Our input consists of two sets of entities. We apply state-of-the-art NLP and IR techniques to automatically infer commonsense relations between the entities using a variety of data sources, and construct a mapping between the domains. Importantly, we do not require identical phrasings of relations. Moreover, our output is interpretable, showing how the mapping was chosen.

Unlike previous works, we drop the strong *bijectivity* assumption and let the algorithm decide which entities to include in the mapping. Our algorithm can also generate new *suggestions* for the non-mapped entities. This paves the road to algorithms that can handle even more limited input – for example, using domain *names* (solar system, atom) as input, or just a single mapped entity pairs (e.g., turn white blood cells into policemen and see how the analogy unfolds). Our contributions are:

- A novel, scalable and interpretable approach for automatically mapping two domains based on commonsense *relational* similarities. Our algorithm handles partial mappings and suggests additional entities.
- We extend the work of Romero and Razniewski (2020) to discover salient knowledge about pairs of entities.
- Our model's accuracy is 81.2% on simple, 2x2 problems. On larger problems, it achieves 77.8% perfect mappings (guess level: 13.1%). In another experiment, we outperform humans (90% vs. 70.2%) and demonstrate that our au-

tomatic suggestions resemble human suggestions. We release code and data¹. 117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

2 Problem Definition

An analogy is a mapping from a base domain \mathcal{B} into a target domain \mathcal{T} . The mapping is based on *relations*, not object attributes. Base objects are not mapped into objects that resemble them; rather, there is a common *relational structure*, and they are mapped to objects that play similar roles. We follow the formulation of Sultan and Shahaf (2022), brought here for completeness:

Entities and Relations. Let $\mathcal{B} = \{b_1, ..., b_n\}$ and $\mathcal{T} = \{t_1, ..., t_m\}$ be two sets of entities. For example: $\mathcal{B} = \{$ sun, Earth, gravity, solar system, Newton $\}, \mathcal{T} = \{$ nucleus, electrons, electricity, atom, Faraday $\}.$

Let \mathcal{R} be a set of relations. A relation is a set of ordered entity pairs. The exact representation is purposely vague, as we do not restrict ourselves to strings, embeddings, etc. Intuitively, relations should capture notions like "revolve around".

In our example, relations between \mathcal{B} and \mathcal{T} include the *Earth* revolve around the *sun*, like *electrons* orbit the *nucleus*; the *Earth* creates a force field of *gravity*, similar to *electrons* creating *electric force* fields; the *sun* and the *Earth* are part of the *solar system*, as the *nucleus* and *electrons* are part of the *atom*; *Newton* discovered *gravity*, as *Faraday* is credited with discovering *electric force*.

Note that relation is an asymmetric function, as the pairs are ordered; e.g., Newton discovered gravity, but gravity did not discover Newton.

Slightly abusing notation, we denote the *set* of relations that hold between two entities e_1, e_2 as

113

114

115

¹shorturl.at/ADIPR

B	Mapping	\mathcal{T}
Sun	\rightarrow	Nucleus
Earth	\rightarrow	Electrons
Gravity	\rightarrow	Electric force
Solar system	\rightarrow	Atom
Newton	\rightarrow	Faraday

Table 1: Illustration of a relational analogy between the solar system and the atom.

 $\mathcal{R}(e_1, e_2) \subseteq 2^{\mathcal{R}}$. For example, $\mathcal{R}(earth, sun)$ contains {revolve around, attracted to}, etc. For clarity, we sometimes use \mathcal{R}_B , \mathcal{R}_T to emphasize that the entities belong to the \mathcal{B} , \mathcal{T} domain.

150

151

152

153

154

155

156

157

159

160

161

162

163

164

165

167

168

169

170

171

172

173

174

175

176

177

178

179

181

182

Similarity. Let sim be a similarity metric between two sets of relations, sim : $2^{\mathcal{R}} \times 2^{\mathcal{R}} \rightarrow [0, \infty)$. Intuitively, when applied to singletons, we want our similarity metric to capture how relations are like each other. For example, "revolve around" is similar to "orbit" and (to a lesser degree) "spiral". When applied to sets of relations, we want sim to be higher if the two sets share many distinct relations. For example, {revolve around, attracted to} should be more similar to {orbit, drawn into} than to {revolve around, orbit} (as the last set does not include any relation similar to attraction). In Section 3.2 we present our sim implementation.

Given one pair from \mathcal{B} and one from \mathcal{T} , we define similarity in terms of their relations. Since \mathcal{R} is asymmetric, we consider both directions:

$$sim^{*}(b_{1}, b_{2}, t_{1}, t_{2}) =$$

$$sim(\mathcal{R}_{B}(b_{1}, b_{2}), \mathcal{R}_{T}(t_{1}, t_{2})) +$$

$$sim(\mathcal{R}_{B}(b_{2}, b_{1}), \mathcal{R}_{T}(t_{2}, t_{1}))$$

Objective. Output a mapping $\mathcal{M} : \mathcal{B} \to \mathcal{T} \cup \bot$ such that no two \mathcal{B} entities mapped to the same \mathcal{T} entity (Table 1). Mapping into \bot means the entity was not mapped to any entity in the \mathcal{T} domain.

We look for the mapping \mathcal{M}^* that captures the best inter-domain analogical structure similarity by maximizing the relational similarity:

$$\arg\max_{\mathcal{M}} \sum_{j=1}^{n-1} \sum_{i=j+1}^{n} sim^{*}(b_{j}, b_{i}, \mathcal{M}(b_{j}), \mathcal{M}(b_{i}))$$

Note: if b_i or b_j maps to \perp , sim^* is defined to be 0.

3 Analogous Matching Algorithm

We wish to find the best \mathcal{B} to \mathcal{T} mapping. We first extract relations between entity pairs from the same

domain (Section 3.1). Then, we compute similarity between entity pairs that could be mapped (Section 3.2). Finally, we build the mapping (Section 3.3). 183

184

185

186

187

188

189

190

191

192

193

194

195

196

197

198

199

200

202

203

204

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

225

226

227

229

230

3.1 Relation Extraction

Automatically extracting relations is a key part of our algorithm. We focus on *commonsense* relations (e.g., the Earth *revolves around* the sun), as opposed to situational relations (e.g., the book is on the table). This broadly falls under open information extraction (OIE), the task of generating a structured representation of the information in a text. There has been a lot of work in this area, especially attempts to automate the construction of commonsense datasets (Etzioni et al., 2008, 2004; Yates et al., 2007; Lenat et al., 1985; Sap et al., 2019).

Given two entities, we automatically extract relations from multiple sources:

ConceptNet. A commonsense dataset, containing about 1.5M nodes (Liu and Singh, 2004). For each entity, we receive a list of (predicate, entity), which we filtered to match the second entity (single or plural form). The predicates serve as our relations.

Open Information Extraction. A database automatically extracted from a large web corpus (Etzioni et al., 2008). It contains over 5B triplets of (subject, predicate, object). We searched for a match between both entities in the (subject, object) fields, and used the predicates as our relations.

GPT-3². We used a generative pretrained large language model (LM) as a knowledge base in a few-shot manner (Petroni et al., 2019; Brown et al., 2020b). We input a prompt of four analogies, e.g., "Q: What is the relation between gravity and Newton?, A: Newton discovered gravity. A: Newton invented gravity." (see Section A.2.3 for full prompt). GPT-3 outputs up to 3 sentences per query. We kept only sentences of the form <entity> <text> <entity>, treating the <text> as the relation.

Quasimodo. A commonsense knowledge base that focuses on salient properties of objects (Romero and Razniewski, 2020). It contains more than 3.5M triplets of (subject, predicate, object). It considers *questions* instead of statements. For instance, if people search for an answer to "why is the sky blue?", this implies that the sky is blue. Whenever our two entities appeared in the (subject, object) fields, we extracted their predicates as relations.

²GPT-3 is the only data source that is not freely available. All queries needed for this paper accumulated to less than \$50.



Figure 2: Quasimodo++. Example regex used to extract suggestions from Google ("<question> <entity1> .* <entity2>"). We use questions such as "why does", "why did" and "how does".

Quasimodo++. A relation extraction method that we develop, inspired by Quasimodo. Quasimodo was constructed using questions about a single entity; we extended it to questions exploring relations between *pairs* of entities. We used Google's query auto-completion to tap into the query logs, asking questions containing both desired entities, such as "How does earth * sun", "How is earth * sun", and "Why does sun * earth" for every pair of entities (see Figure 2 for an example). The exact regular expressions we used can be found in Section A.1.

Our algorithm is easy to extend to new sources. We expect robustness will increase with coverage.

3.2 Scoring Entity Pairs

235

239

240

241

243

244

245

247

249

252

253

257

261

263

264

265

266

We wish to calculate $sim^*(b_i, b_j, t_k, t_p)$ for $b_{i,j} \in \mathcal{B}$, $t_{k,p} \in \mathcal{T}$, $1 \le i < j \le n$, $1 \le k \ne p \le m$.

In Section 2 we specified desiderata of *sim*, especially that it is higher if the two sets share many distinct relations. We turn to our implementation.

Without loss of generality, let us consider $sim(\mathcal{R}_B(b_1, b_2), \mathcal{R}_T(t_1, t_2))$. We first extract all relations $\mathcal{R}_B(b_1, b_2), \mathcal{R}_T(t_1, t_2)$. Next, we calculate the score between each relation in $\mathcal{R}_B(b_1, b_2)$ and each relation in $\mathcal{R}_T(t_1, t_2)$. We create a complete bipartite graph where the left side nodes are the relations of $\mathcal{R}_B(b_1, b_2)$, and the right side nodes are the relations of $\mathcal{R}_T(t_1, t_2)$ (Figure 3). The edge weights (w) are the cosine similarity of the nodes' sBERT embedding (Reimers and Gurevych, 2019).

We remove non-informative relations by extracting the top-frequent n-grams $(n = \{1, 2, 3, 4\})$ from Wikipedia and setting their score to zero. Edges that did not reach a threshold (chosen using fine-tuning, see Section 3.3) were set to zero.

Next, we cluster similar relations on each side (e.g., "revolve around" and "circle around") to avoid double-counting. We use hierarchical agglomerative clustering based on the cosine embed-



Figure 3: Left: partial relations of *Earth:sun*. Right: partial relations of *electron:nucleus*. This is the result of the maximum weighted match on the clusters. Colors correspond to clusters.

ding similarity (threshold = 0.5; see Section 3.3). The weight of edges between two clusters is the maximal weight of an edge between their nodes (see Figure 3; colors correspond to clusters).

Finally, we apply Maximum-Weight Bipartite Matching on the clusters (see Section 3.3). The similarity score $sim(\mathcal{R}_B(b_1, b_2), \mathcal{R}_T(t_1, t_2))$ is defined as the sum of the remaining edges.

3.3 Building a Mapping

Using the score mappings between pairs, we can compose larger mappings. We use beam-search, starting from the most promising pair-mappings of Section 3.2. In each iteration, we expand the 20 most promising partial mappings, testing each possible mapping between single entities of \mathcal{B} and \mathcal{T} (that are consistent with the current partial mapping – i.e., a \mathcal{B} entity cannot map to multiple \mathcal{T} entities). When expansions stop increasing the score, we stop the search and select the highest score mapping.

Figure 4 shows a snippet from our UI. Input appears on the top. FAME's output mapping is represented as a graph: nodes correspond to single entity mappings (e.g., sun to nucleus). Edges represent the shared relational structure. Each edge contains some of the shared relations between the mapped pairs corresponding to its endpoints (e.g., "more massive than") and their similarity score (note the edges are directional). To ease visualization, we show at most two relations per edge. The weight of an edge corresponds to its strength.

A note on the solution space. In other works n = m and \mathcal{M} is a *bijective* function, and the solution space's cardinality is n!. We allow for



Figure 4: A snippet from our UI. Top: Input. Bottom: The mapping our algorithm found (output), is represented as a graph. Nodes correspond to mappings between single entities (e.g., sun to nucleus). Each edge is annotated with some of the shared relations between the mapped pairs corresponding to its endpoints and their similarity score. For the sake of visualization, we show at most two relations for each edge. Edge weight corresponds to strength.

 $n \neq m$ and unmapped entities. Without loss of generality let $n \leq m$. The cardinality is then $\left(\sum_{i=0}^{n} {n \choose i} \frac{m!}{(m-i)!}\right) - (n \cdot m)$, where *i* is the number of matched entities. We subtract $n \cdot m$ because we do not allow for a mapping of size 1; our algorithm starts by mapping pairs and then adds single-entity mapping at each iteration of the beam search.

We note that relaxing the bijective constraint drastically increases the space. For n = 7, n! = 5,040, while our space is of size 130,922.

Fine-tuning. We constructed a new dataset to fine-tune our model's hyper-parameters (See Appendix A). The dataset contains 36 analogical mapping problems created by ten volunteers, not from our research team. We showed them example analogies and asked them to generate new ones. An expert from our team verified their output, discarding 4 analogies. Domain size was 3-5 (average 3.4).

On the problems generated by the volunteers, we achieve 83.3% perfect mappings (whole mapping is correct). If we consider single mappings separately, the algorithm achieves 89.4% accuracy.

4 Entity Suggestion

One of the main limitations of previous analogical mapping algorithms is their inability to automatically expand analogies. This is especially interesting in our case, as we allow for unmapped entities; suggesting new entities could identify potential mapping candidates for the unmapped entities.

For example, let $\mathcal{B} = \{$ sun, Earth, gravity, Newton $\}$ and $\mathcal{T} = \{$ nucleus, electron, electricity $\}$. The correct mapping is sun \rightarrow nucleus, Earth \rightarrow electron, gravity \rightarrow electricity, leaving Newton with no

mapping. Our goal is to suggest candidate entities that preserve the relational structure.

Intuitively, we look at the relations Newton shares with other \mathcal{B} entities (e.g., discovered gravity), and try to see which \mathcal{T} entity plays a corresponding role (i.e., who discovered electricity?).

337

338

339

341

342

343

344

346

347

348

351

352

353

354

355

356

357

358

359

360

362

363

364

366

More formally, suppose we wish to find candidates t^* for mapping to b_n . We first extract the relations of $R_b(b_i, b_n)$, $\forall i \in [n]$ (denoted as R_{b_i}). We then iterate over all relations $r \in R_{b_i}$ and use the pair $\{\mathcal{M}(b_i), r\}$ to extract suggestions for t^* .

We use Open Information Extraction, Quasimodo, and Quasimodo++. While our method was previously used to extract *relations* given a pair of two entities, we now use it to extract *entities* given a pair of {entity, relation}. This entails filtering on the predicate field in our commonsense datasets and changing the queries in Quasimodo++.

As suggestions tend to be noisy, we cluster all extracted entities (similarly to the clustering from Section 3.2). We remove clusters of size < 2.

For each suggestion cluster, we rerun our analogous matching algorithm with a representative entity from that cluster (the closest to the cluster's center of mass). We pick the cluster whose representative resulted in the mapping with the highest score. As the commonsense datasets we work with operate mostly on string matching, small changes (e.g., Benjamin Franklin/Ben Franklin) could sometimes result in slightly different results. Thus, we perform one final round, with *all* entities from our chosen cluster, and pick the highest score mapping.

Sources	Near	Far	Extended
All	85%	77.5%	$\mathbf{77.8\%}$
All-ConceptNet	85%	77.5%	$\mathbf{77.8\%}$
All-Open IE	85%	67.5%	58.3%
All-Quasimodo	85%	$\mathbf{77.5\%}$	72.2%
All-Quasimodo++	80%	72.5%	72.2%
All-GPT3	57.5%	50%	66.7%

Table 2: Ablation study on the 2x2 near and far problems and our extended set, leaving out knowledge sources. Results show the importance of the generative LM approach (GPT-3) as a knowledge source. Open Information Extraction also contributes much, especially for the complex analogies (2x2-far and extended).

5 Evaluation

In this section, we evaluate FAME. We test its ability to identify the correct mapping (Section 5.1), and compared it to related works (Section 5.2) and human performance (Section 5.3).

5.1 Performance on Analogy Problems

2x2 problems. One of the things that might have held computational analogy back is the lack of high-quality, large-scale datasets. Most datasets are small and focus on classical 2x2 problems (A : B :: C : D), similar to SAT questions.

We start by testing FAME on this standard type of analogies. We use 80 problems from Green et al. (2010), split into 40 near and 40 far analogies (e.g., for "answer:riddle", near analogy is "solution:problem", far analogy is "key:lock"). While the dataset is small, we believe it is still interesting to explore. Our algorithm managed to perfectly map 85% of near analogies and 77.5% of far ones. Random guess baseline is 33.3% (Section 3.3).

Extended problems. Encouraged by the results of the 2x2 problems, we explore more complex problems. We decided to extend the Green far analogies (which are harder than the near ones). We had three experts go over the dataset together and brainstorm potential extensions. On four problems, the experts did not manage to agree on any additional mappings, leaving us with 36 extended problems (average domain size 3.3).

Our algorithm perfectly mapped 77.8% of the extended problems. Random baseline is 13.1% on average. As we relax the bijection assumption, our algorithm's average guess level is 2.2% (Section 3.3). If we look beyond the top-rated solution, our algorithm has the correct solution in its top-2

guesses 83.3% of the time and 91.7% for top-3.

Error analysis. We found 3 main causes of error:

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

- Coverage (for example, we could not find a relation between "hoof" and "hoofprint"). This prompted us to ablate the knowledge sources FAME uses (Table 2). Results show the importance of the generative LM approach. Open IE is also important, especially for the more complex analogies (far and extended). Some sources, such as ConceptNet, did not seem to contribute much.
- **Noisy relations** that are either peculiar or plain wrong (e.g., "a footballer can iron").
- Embedding similarity (for example, "produce" and "is produced by" have a high similarity score). This is exacerbated by **ambigu**ity (e.g., the word "pen" referred to "pigpen" and not to the writing instrument).

5.2 Comparison to Related Work

SME line of work. We had difficulty comparing FAME to SME (Falkenhainer et al., 1989) and its extensions, due to their complex input requirements. LRME (Turney, 2008) is closest to our setting, but no code or demo is available. Thus, we compare to their published results on a set of 20 problems.

LRME's entities include nouns, verbs, and adjectives. Since FAME expects noun phrases, we filtered out all other input terms (one problem has only a single noun, so we are left with 19 problems). It is hard to compare in this setup (and unfortunately, authors did not report which partial mappings were correct). Still, LRME's accuracy was 75%, whereas FAME achieved 84.2%.

While the size of the problems is smaller when restricted to nouns, we believe the noun-only setting is harder. The verbs and adjectives often provide hints that significantly constrain the search space. For example, in problem A6 (Turney, 2008) (mapping a projectile to a planet) there is one adjective in each domain (parabolic, elliptical). Those adjectives can only apply to one or two of the nouns (i.e., you cannot have parabolic earth, air, or gravity), effectively giving away the noun mapping.

As a side note, we also believe that our nounonly input is a cleaner problem setting, as it is often easier to automatically identify the entities in a domain than to identify the attributes and verbs relevant for the analogy. In the words of LRME's authors, "LRME is not immune to the criticism of Chalmers et al. (1992), that the human who

375

378

381

394

399

400

401

367

Algorithm	Near	Far	Extended
FAME	85.0%	77.5%	77.8%
GPT-3 ":"	92%	80 %	44%
GPT-3 "->"	88%	80 %	58%

Table 3: Comparison of FAME and GPT-3. GPT-3 does well on the 2x2 datasets (far and near). We note that data leakage is a concern. GPT-3's performance sharply drops on the extended problem set, where problems are bigger and do not appear on the web.

generates the input is doing more work than the computer that makes the mapping." We believe FAME is a step in the right direction in this regard. **Pretrained LMs.** In the absence of a baseline, we turn to a generative pretrained large LM known to have impressive commonsense abilities – GPT-3. We used 4 random examples from the fine-tuning dataset. After some experimentation with prompt engineering, we chose two variants (see A.2.3).

The results are summarized in Table 3. GPT-3 does well on the 2x2 datasets (Green et al., 2010). However, both datasets appear on the web, and perhaps GPT-3 was exposed to them during training (data leakage). In particular, we found some of the answers via a simple web search (Figure A.6).

Moreover, GPT-3's performance drops on the extended set, where problems are complex and do not appear on the web. Interestingly, it does not even manage to return a valid mapping in some of the cases. This exercise improves our understanding of FAME's strengths and weaknesses.

E-KAR dataset. Chen et al. (2022) recently released a relevant dataset, E-KAR, for rationalizing analogical reasoning. The dataset consists of multiple-choice problems from civil service exams in China. For example, for the source triplet "tea:teapot:teacup", the correct answer is "talents:school:enterprise". The reasoning is that both teapot and teacup are containers for tea. After the tea is brewed in the teapot, it is transported into the teacup. Similarly, both school and enterprise are organizations. After talents are educated in school, they are transported into enterprise³.

The E-KAR test set has no labels, so we used their validation set (N=119) to test FAME. As our task is different, we only took source entities (as \mathcal{B}) and entities from the correct answer (as \mathcal{T}). We filtered questions without nouns, resulting in N=101. FAME found the right mapping 68.3% of the time. A closer examination of FAME's mistakes revealed that ~ 75% of them occurred due to relation *types* that are not at all covered by our framework: either ternary relations (soldier:doctor:military doctor \rightarrow car:electric vehicle:electric car; the last term is a combination of the first two) or relations based on sharing some attribute (so "both containers for holding tea" is mapped to "both are organizations"). Some of the attribute-based mappings work at the whole-set level, so each entity on \mathcal{B} could map to each entity on \mathcal{T} (yellow:red:white \rightarrow sad:happy:angry). Thus, we conclude there is a big gap between FAME and E-KAR's assumptions.

5.3 Comparison to People

We compare FAME with *human thinking* in a 2phase experiment⁴. In the closed-world phase, the participants received ten structure mapping problems, in which they were asked to match instances from \mathcal{B} to \mathcal{T} . The domains included between 3-5 entities (Table A.4). Participants were instructed to map each \mathcal{B} entity into exactly one \mathcal{T} entity.

In the open-world phase, participants received five *mapped* problems, but one entity was left blank (Table A.5). Participants were instructed to fill in the blank with an entity that preserves the analogy.

Participants. We recruited 304 participants using social media. The compensation was a chance to win one of three \$30 vouchers. 76.6% of our participants were between the ages 18-35 and 17.2% are between 36-45 (self-reported).

Closed-world mapping. FAME missclassified one problem compared to gold standard (A9, Table A.4), achieving 90% accuracy (human baseline was 70.2%; see full distribution in Table A.4).

Problem A6 has the lowest human accuracy (35.5%), and is also the largest one $(|\mathcal{B}| = |\mathcal{T}| = 5)$. A closer examination of its confusion matrix reveals that while FAME correctly mapped *water* to *heat* and *pressure* to *temperature*, 15% of people switched the two. This might be due to the strong semantic pairing of *water* and *temperature*. FAME is immune to this, as it relays on *relations*.

On average, each participant mapped the problem the same as FAME 78% of the times. Overall, FAME outperforms humans, and most of the disagreement is due to human errors.

Open-world: entity suggestion. We presented

³Interestingly, the authors of this paper thought that the "passengers:bus:taxi" answer was the correct one, based on containment and size relations.

⁴The experiment received ethics committee approval. See full instructions in Section A.4.



Figure 5: Word cloud of human completions for *B*1 (Table A.6). While most responses were from the same semantic domain, some were creative and appropriate (e.g., treasure chest, jewelry box, car).

participants with five *mapped* problems where one entity was left blank (Table A.5) and asked them to fill in the black while preserving the analogy.

For all five problems, an entity from FAME's top two completions appeared in humans' top three completions (Table A.6). Meaning, our algorithm's top suggestions are similar to humans'. Only in one example (*B*5) one of the top two algorithm's completions appeared *third* in humans' (in the rest it is first or second). We suspect that *gravity* and *Newton* reminded participants of the term *apple*.

Figure 5 shows a word cloud for answers to problem B1. While most responses are quite similar, some participants returned creative and appropriate solutions (e.g., treasure chest, jewelry box, car).

6 Related Work

Computational analogy-making dates back to the 1960s (Evans, 1964; Reitman, 1965). Analogymaking approaches are broadly categorized as symbolic, connectionist, and hybrid (French, 2002; Mitchell, 2021; Gentner and Forbus, 2011).

Symbolic approaches usually represent input as structured sets of logic statements. Our work falls under this branch, as well as SME (Falkenhainer et al., 1989) and its follow-up work. LRME (Turney, 2008) is the closest to our work, as it automatically extracts the relations. Unlike FAME, LRME requires exact matches of relations across different domains. We also focus on nouns only, making the problem harder, and relax the bijection assumption, allowing for automatically extending analogies.

NLP. Analogy-making received relatively little attention in NLP. The best-known task is word analogies, often used to measure embeddings' quality (inspired by Word2Vec's "*king - man + woman = queen*" example (Mikolov et al., 2013)). Follow-up work explored embeddings' linear algebraic structure (Arora et al., 2016; Gittens et al., 2017; Allen and Hospedales, 2019) or compositional nature (Chiang et al., 2020), neglecting relational similarity. A recent work on analogies between procedural texts (Sultan and Shahaf, 2022) did study relational similarity, but extracted the relations from the input texts, with no commonsense augmentations. 574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

597

598

599

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

Recently, there have been efforts to study LMs' analogical capabilities (Ushio et al., 2021; Brown et al., 2020a). Findings indicate they struggle with abstract and complex relations and results depend strongly on LM's architecture and parameters.

Kittur et al. (2019) combined NLP and crowds for product analogies without explicitly modeling entities and relations, but instead automatically extracting *schemas* of the product.

7 Conclusions and Future Work

Detecting deep structural similarity across distant domains and transferring ideas between them is central to human thinking. We presented FAME, a novel method for analogy making. Compared to previous works, FAME is more expressive, scalable, robust and interpretable. It also allows partial matches and automatic entity suggestions to extend the analogies.

FAME correctly maps 81.2% of classical 2x2 analogy problems. On larger problems, it achieves 77.8% perfect mappings (mean guess level: 13.1%). FAME also outperforms humans in solving analogy mapping problems (90% vs. 70.2%). Interestingly, our automatic suggestions of new entities resemble those suggested by humans.

In future work, we plan to improve coverage and extend our framework to more than just binary relations, as sometimes the key to an analogy is a relation involving more than two objects. In addition, we plan to improve our similarity measure, to address both context (to solve ambiguity) and the difference between active and passive relations. We plan to explore different forms of input, such as algorithms that take as input very partial domains, perhaps even just domain *names* (solar system, atom) and populate the domains with entities, or algorithms incorporating *user feedback*.

To conclude, we hope FAME will pave the way for analogy-making algorithms that require lessrestrictive inputs and can scale up and tap into the vast amount of potential inspiration the web offers, augmenting human creativity.

8

573

538

8 Ethical Considerations & Limitations

While FAME can assist humans by inspiring nontrivial solutions to problems, it might also be somewhat misleading. It has been shown that humans struggle with detecting caveats in presented analogies (Holyoak et al., 1995). For example, the cardiovascular system is often taught to medical students in terms of water supply system (Swain, 2000). However, this analogy might also confuse them, as it ignores important differences between water and blood (e.g., blood clots). Thus, while our output is interpretable, it might still mislead people, and it is important to alert the users to this possibility.

627

639

641

644

647

Another ethical consideration is the fact that FAME's coverage highly depends on external resources (ConceptNet, Google AutoComplete, etc.). This might be particularly problematic when applied to low-resource languages. As the relations we look for are commonsense relations, rather than cultural or situational ones, using automatic translation might ameliorate the problem.

Lastly, we also note these resources evolve over time, and thus if one is interested in reproducibility, it is necessary to save the extracted relations.

References

Carl Allen and Timothy Hospedales. 2019. Analogies explained: Towards understanding word embeddings. In *Proceedings of the 36th International Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 223– 231. PMLR. 649

650

651

652

653

654

655

656

657

658

659

660

661

662

663

664

665

666

667

668

669

670

671

672

673

674

675

676

677

678

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

703

- Sanjeev Arora, Yuanzhi Li, Yingyu Liang, Tengyu Ma, and Andrej Risteski. 2016. A latent variable model approach to pmi-based word embeddings. *Transactions of the Association for Computational Linguistics*, 4:385–399.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020a. Language models are few-shot learners. In Advances in Neural Information Processing Systems, volume 33, pages 1877–1901. Curran Associates, Inc.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020b. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- David J Chalmers, Robert M French, and Douglas R Hofstadter. 1992. High-level perception, representation, and analogy: A critique of artificial intelligence methodology. *Journal of Experimental & Theoretical Artificial Intelligence*, 4(3):185–211.
- Jiangjie Chen, Rui Xu, Ziquan Fu, Wei Shi, Zhongqiao Li, Xinbo Zhang, Changzhi Sun, Lei Li, Yanghua Xiao, and Hao Zhou. 2022. E-kar: A benchmark for rationalizing natural language analogical reasoning. *arXiv preprint arXiv:2203.08480*.
- Hsiao-Yu Chiang, Jose Camacho-Collados, and Zachary Pardos. 2020. Understanding the source of semantic regularities in word embeddings. In *Proceedings of the 24th Conference on Computational Natural Language Learning*, pages 119–131, Online. Association for Computational Linguistics.
- Morteza Dehghani and Ken Forbus. 2009. Qcm: A qp-based concept map system. In *the 23nd International Workshop on Qualitative Reasoning (QR09)*, pages 16–21.
- Oren Etzioni, Michele Banko, Stephen Soderland, and Daniel S Weld. 2008. Open information extraction from the web. *Communications of the ACM*, 51(12):68–74.

705

758

pages 100-110. Thomas G Evans. 1964. A heuristic program to solve geometric-analogy problems. In Proceedings of the April 21-23, 1964, spring joint computer conference, pages 327-338.

Oren Etzioni, Michael Cafarella, Doug Downey, Stan-

ley Kok, Ana-Maria Popescu, Tal Shaked, Stephen

Soderland, Daniel S Weld, and Alexander Yates.

2004. Web-scale information extraction in know-

itall: (preliminary results). In Proceedings of the

13th international conference on World Wide Web,

- Brian Falkenhainer, Kenneth D Forbus, and Dedre Gentner. 1989. The structure-mapping engine: Algorithm and examples. Artificial intelligence, 41(1):1-63.
- Kenneth Forbus, Jeffrey Usher, Andrew Lovett, Kate Lockwood, and Jon Wetzel. 2011. Cogsketch: Sketch understanding for cognitive science research and for education. Topics in Cognitive Science, 3(4):648-666.
- Robert M French. 2002. The computational modeling of analogy-making. Trends in cognitive Sciences, 6(5):200-205.
- Dedre Gentner. 1983. Structure-mapping: A theoretical framework for analogy. Cognitive science, 7(2):155-170.
- Dedre Gentner and Kenneth D Forbus. 2011. Computational models of analogy. Wiley interdisciplinary reviews: cognitive science, 2(3):266–276.
- Alex Gittens, Dimitris Achlioptas, and Michael W Mahoney. 2017. Skip-gram- zipf+ uniform= vector additivity. In Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 69-76.
- Adam E Green, David JM Kraemer, Jonathan A Fugelsang, Jeremy R Gray, and Kevin N Dunbar. 2010. Connecting long distance: semantic distance in analogical reasoning modulates frontopolar cortex activity. Cerebral cortex, 20(1):70-76.
- Douglas R Hofstadter and Emmanuel Sander. 2013. Surfaces and essences: Analogy as the fuel and fire of thinking. Basic books.
- Keith J Holyoak. 1984. Analogical thinking and human intelligence. Advances in the psychology of human intelligence, 2:199-230.
- Keith J Holyoak, Paul Thagard, and Stuart Sutherland. 1995. Mental leaps: analogy in creative thought. Nature, 373(6515):572-572.
- Aniket Kittur, Lixiu Yu, Tom Hope, Joel Chan, Hila Lifshitz-Assaf, Karni Gilon, Felicia Ng, Robert E Kraut, and Dafna Shahaf. 2019. Scaling up analogical innovation with crowds and ai. Proceedings of the National Academy of Sciences, 116(6):1870-1877.

Douglas B Lenat, Mayank Prakash, and Mary Shepherd. 1985. Cyc: Using common sense knowledge to overcome brittleness and knowledge acquisition bottlenecks. AI magazine, 6(4):65-65.

759

760

761

763

764

766

767

768

769

770

771

772

773

774

775

776

779

780

781

782

783

784

785

786

787

788

791

792

793

794

795

796

797

798

799

800

801

802

803

804

805

806

807

808

809

810

- Hugo Liu and Push Singh. 2004. Conceptnet-a practical commonsense reasoning tool-kit. BT technology journal, 22(4):211–226.
- Donald G McNeil Jr and Mr Odón. 2013. Car mechanic dreams up a tool to ease births. The New York Times, 13.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. Advances in neural information processing systems, 26.
- Marvin Minsky. 1988. Society of mind. Simon and Schuster.
- Melanie Mitchell. 2021. Abstraction and analogymaking in artificial intelligence. Annals of the New York Academy of Sciences, 1505(1):79–101.
- Fabio Petroni, Tim Rocktäschel, Patrick Lewis, Anton Bakhtin, Yuxiang Wu, Alexander H Miller, and Sebastian Riedel. 2019. Language models as knowledge bases? arXiv preprint arXiv:1909.01066.
- PJM. 1966. Models and analogies in science.
- Nils Reimers and Iryna Gurevych. 2019. Sentencebert: Sentence embeddings using siamese bertnetworks. arXiv preprint arXiv:1908.10084.
- Walter Ralph Reitman. 1965. Cognition and thought: an information processing approach.
- Julien Romero and Simon Razniewski. 2020. Inside quasimodo: Exploring construction and usage of commonsense knowledge. In Proceedings of the 29th ACM International Conference on Information & Knowledge Management, pages 3445-3448.
- Maarten Sap, Ronan Le Bras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A Smith, and Yejin Choi. 2019. Atomic: An atlas of machine commonsense for ifthen reasoning. In Proceedings of the AAAI Conference on Artificial Intelligence, volume 33, pages 3027-3035.
- Oren Sultan and Dafna Shahaf. 2022. Life is a circus and we are the clowns: Automatically finding analogies between situations and processes. Proceedings of the 2022 conference on empirical methods in natural language processing (EMNLP).
- David P Swain. 2000. The water-tower analogy of the cardiovascular system. Advances in Physiology Education, 24(1):43-50.
- Peter D Turney. 2008. The latent relation mapping engine: Algorithm and experiments. Journal of Artificial Intelligence Research, 33:615–655.

Asahi Ushio, Luis Espinosa Anke, Steven Schockaert, and Jose Camacho-Collados. 2021. Bert is to nlp what alexnet is to cv: Can pre-trained language models identify analogies? In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 3609–3624.

812

813

814 815

816

819

820

821

822 823

824

825

826

Alexander Yates, Michele Banko, Matthew Broadhead, Michael J Cafarella, Oren Etzioni, and Stephen Soderland. 2007. Textrunner: open information extraction on the web. In Proceedings of Human Language Technologies: The Annual Conference of the North American Chapter of the Association for Computational Linguistics (NAACL-HLT), pages 25–26.

A Implementation Details

We fine-tune our model using 36 problems described in Section 3.3.

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

We used the pre-trained model *msmarco-distilbert-base-v4* which is based on sBERT (Reimers and Gurevych, 2019). We set the similarity threshold (the similarity between two relations) to be **0.2** (range checked: 0-0.6). We set the number of top n-grams which was filtered (the top frequencies n-grams in Wikipedia) to **500**. The clustering distance threshold is set to **0.5** (range checked: 0.3-0.9). The number of clusters we consider when computing the sum is set to **3** (range checked: 1-maximum number of clusters). We set the beam search size to **20** (range checked: 1-40). All of these parameters describes in Section 3.

We provide access to our anonymous repository can be found¹. We note that the usage of Docker is not supported in this version for the purpose of maintaining anonymity. However, the algorithmic content is available.

A.1 Quasimodo++ regular expressions

We use the following regex for our Quasimodo++: "<question> <prefix> <entity1> .* <entity2>". The questions we used are: {"why do", "why is", "why does", "why does it", "why did", "how do", "how is", "how does", "how does it", "how did"}. The prefix is optional and can be {"a", "an" and "the"}. We use both singular and plural forms of the entities.

A.2 GPT-3

A.2.1 Prompts used for relation extraction

The prompt used for GPT-3 is: 859 Q: What are the relations between a blizzard and 860 snowflake? 861 A: A blizzard produces snowflakes. 862 A: A blizzard contains a lot of snowflakes. 863 864 O: What are the relations between an umbrella and 865 rain? 866 A: An umbrella protects from rain. 867 A: An umbrella provides adequate protection from 868 rain 869 870 Q: What are the relations between a movie and 871 screen? 872 A: A movie displayed on a screen. 873 A: A movie can be shown on a screen. 874 875





Figure 6: Looking for analogies from the original Green eval dataset online.

876	Q: What are the relations between Newton and
877	gravity?
878	A: Newton discovered gravity.
879	A: Newton invented gravity.
880	
881	Q: What are the relations between an electron and
882	nucleus?
883	A: An electron revolves around the nucleus.
884	A: An electron is much smaller than the nucleus.
885	A: An electron attracts the nucleus.

Q: What are the relations between water and a pipe?

A: Water flows through the pipe.

A: Water passes through the pipe.

886

892

900

901

902

903

A.2.2 Prompts used for baseline comparison

After some experimentation with prompt engineering, we chose two variants of the prompt:

Q: Find an analogical mapping between the entities "eraser", "paper" and "pencil" and the entities "keyboard", "delete" and "screen".

A: eraser:pencil:paper::delete:keyboard:screen

(or)

A: eraser -> delete, pencil -> keyboard, paper -> screen

A.2.3 Possible leakage

Example answers for chosen analogies from Green eval dataset found via a simple web search can be found in Figure 6

A.3 Repository

905To ease the access and usage of our code we use906Docker. Its main goal is to shift the cross-platform907installation burden from the user to the developer.908Unfortunately, we cannot share our Docker due to

anonymity concerns (username). We will include it in the non-anonymized version.

We provide a React based web interface, currently available only locally. This system is used to visualize the graphs created by the algorithm's mapping output. In addition, it visualizes the relations between entities, their similarity, and the clustering. This interface is useful for assisting in developing, debugging and understanding the algorithm's output. The demo is accessible using our repository¹.

A.4 Experiments

Snippets of the experimental setup (including instructions) can be found in Figures 7, 8.

Table 4 depicts the ten analogical proportion problems used in the *structure mapping* experiment (closed-world mappings in Section 5.2). Accuracy denotes the percentage of human participants who mapped from \mathcal{B} to \mathcal{T} correctly. Results show this task is non-trivial even for humans.

Table 6 illustrates the experimental setup for the second phase of our experiment, in which participants received a solved mapping problem with one entity left out (open-World in Section 5.2).

Table 5 contains all solved analogy problems used in the second phase of the experiment (entity suggestion, see open-World in Section 5.2). Participants were given with the complete mapping, but with a missing entity (as presented here).

A.5 E-kar

Table 7 shows an example of a problematic problemfrom E-KAR dataset.

Instructions			
In this study we are interest For example, we want to ma	sted in analogies mapping between two domains. ap between the car domain, containing car, road, and engine and the boat domain, containing boat, river and sail .		
A possible answer is:			
car road engine	→ boat → river → sail		
The car travels on the road The engine gives the car p	d as the boat sails in the river, so the car maps to the boat and the road maps to the river. power to travel on the road as the sail gives the boat power to sail on the river, so we can map engine to sail.		
In the following we will show you 10 similar analogical mapping problems, your task is to map entities from one side to the other. You should use each entity exactly once, meaning that each entity on the left maps to exactly one entity on the right. There is not necessarily a correct answer, just answer what makes the most sense to you. Feel free to use a dictionary, wikipedia, or any other resources.			
Problem 1			
baker	→		
cake	→		
recipe	→		
ingredients	→		
	(research, scientist, discovery, data)		

Figure 7: Closed-World Mapping: Experiment instructions with the first question.

Instructions				
This part is very similar to th Your task is to fill out the mis	e first par ssing entit	rt, but this time we give you most of the analogical mapping, but leave out one entity. ity. That is, you need to think of a way to complete the given mapping. For example:		
car road engine The mapping expresses an a The car travels on the road We leave the mapping for en We note that you should con If you looked at the engine s You are welcome to be creat	→ b → r → ? nalogy be as the bo gine emp sider all r eparately ive, there	boat river ? etween car and boat : oat sails in the river , so the car maps to the boat and the road maps to the river . pty. What is the equivalent of a car engine in boats? There are multiple things that give boats power, including sail and even engine. relations between car, road and engine. <i>y</i> , you might think about electricity or gears, but that is not the intention. e is no right or wrong here.		
ree nee to use a dictionary,				
Problem 1				
stylist	→ I	landscaper		
hair	→ I	lawn		
gel	→ [

Figure 8: Open-World Entity Suggestion: Experiment instructions with the first question.

		Human Accuracy		
	\mathcal{B}	Mapping	\mathcal{T}	(Guess Level)
	Baker		Scientist	
A 1	Cake	,	Discovery	79.6%
AI	Recipe	\rightarrow	Research	(4.2%)
	Ingredients		Data	
	Eraser		Amnesia	71 707
A2	Pencil	\rightarrow	Memory	(1.7%) (16.7%)
	Paper		Mind	(10.770)
	Jacket		Wound	CO 007
A3	Zipper	\rightarrow	Suture	08.8% (16.7%)
	Cold		Infection	(10.770)
	Train		Signal	74.007
A4	Track	\rightarrow	Wire	(16, 7%)
	Steel		Copper	(10.770)
	Thoughts		Astronaut	F 2.007
A5	Brain	\rightarrow	Space	53.9% (16.7%)
	Neurons		Stars	(10.770)
	Water		Heat	
A6	Pressure		Temperature	or r07
	Bucket	\rightarrow	Kettle	35.5%
	Pipe		Iron	(0.070)
	Rain		Sun	
	Waves		Sounds	
A 77	Water	,	Air	65.1%
A/	Shore	\rightarrow	Ear	(4.2%)
	Breakwater	•	Earplugs	
	Goal		Basket	
	Soccer	,	Basketball	94.1%
Аð	Grass	\rightarrow	Hardwood	(4.2%)
	Feet		Hands	
	Seeds		Ideas	64 507
A9	Fruit	\rightarrow	Product	04.0% (16.7%)
	Bloom		Success	(10.170)
	Morning		Evening	
A 10	Breakfast	,	Dinner	95.1%
AIU	Start \rightarrow	End	(4.2%)	
	Coffee		Wine	

Table 4: The ten analogical proportion problems used in the *structure mapping* experiment. Accuracy denotes the percentage of human participants who mapped from \mathcal{B} to \mathcal{T} correctly. Note that each row under the \mathcal{B} column is mapped to its \mathcal{T} column. Problem's guess level appears in brackets below the accuracy. Results show this task is non-trivial even for humans.

B	Mapping	\mathcal{T}
Electrons	\rightarrow	Earth
Electricity	\rightarrow	Gravity
Faraday	\rightarrow	Newton
Nucleus	\rightarrow	?

Table 5: Solved mapping problem with one missing \mathcal{T} entity. Participants instructed to fill in the missing entity.

	\mathcal{B}	\mathcal{T}	Algorithm	Humans
-	Answer	Key		
R1	Logic	Mechanism		
DI	Riddle	?		
			Problem	Lock (58.9%)
			Lock	Door (11.8%)
			Feedback	Question (4.6%)
	Earth	Electrons		
	Gravity	Electricity		
B2	Newton \rightarrow	Faraday		
	?	Nucleus		
			Sun	Earth's core (15.8%)
			Moon	Apple (13.2%)
			Mars	Sun (10.2%)
	Stylist	Landscaper		
DA	Hair	Lawn		
B3	$Gel \rightarrow$?		
			Fertilizer	Fertilizer (29.3%)
			Water	Lawn Mower (21.1%)
			Lime	Shears (10.2%)
	Chef	Baker		
	Meal	Cake		
B4	Pan \rightarrow	Oven		
	Salt	?		
			Butter	Sugar (63.5%)
			Sugar	Flour (6.9%)
			Onion	Pepper (3.3%)
	Sun	Rain		
P 5	Summer	Winter		
DO	Sunscreen \rightarrow	?		
			Umbrella	Umbrella (51.0%)
			Birds	Coat (20.7%)
			Flooding	Cream (9.9%)

Table 6: Examples used in the second phase of the experiment. Participants were given with the complete mapping, but with a missing entity (as presented here). The algorithm top three completions are sorted according to certainty. Humans' top three completions are sorted according to their frequency in the experiment (in brackets).

	\mathcal{B}	Mapping	\mathcal{T}
Ice		\rightarrow	Grass
Fog		\rightarrow	Tree

Table 7: "ice" and "fog" are different forms of the same substance, and both "ice" and "fog" are natural objects.". "grass" and "tree" are both plants, and "grass" and "tree" are both natural objects.