Causal Parrots: Large Language Models May Talk Causality But Are Not Causal

MORITZ WILLIG^{*} and MATEJ ZEČEVIĆ^{*}, Technical University of Darmstadt, Germany

DEVENDRA SINGH DHAMI, Technical University of Darmstadt and hessian.AI, Germany

KRISTIAN KERSTING, Technical University of Darmstadt, hessian.AI and DFKI, Germany

Some argue scale is all what is needed to achieve AI, covering even causal models. We make it clear that large language models (LLMs) cannot be causal and give reason onto why sometimes we might feel otherwise. To this end, we define and exemplify a new subgroup of Structural Causal Model (SCM) that we call meta SCM which encode causal facts about other SCM within their variables. We conjecture that in the cases were LLM succeed in doing causal inference, underlying was a respective meta SCM that exposed correlations between causal facts in natural language on whose data the LLM was ultimately trained. If our hypothesis holds true, then this would imply that LLMs are like parrots in that they simply recite the causal knowledge embedded in the data. Our empirical analysis provides favoring evidence that current LLMs are even weak 'causal parrots.'

Additional Key Words and Phrases: large language models, correlation of causal facts conjecture, empirical analysis

ACM Reference Format:

Moritz Willig, Matej Zečević, Devendra Singh Dhami, and Kristian Kersting. 2023. Causal Parrots: Large Language Models May Talk Causality But Are Not Causal. 1, 1 (April 2023), 22 pages. https://doi.org/10.1145/nnnnnnnnnnnn

1 INTRODUCTION

Speaking of causality, the Pearlian counterfactual theory of causation has recently found prominent support in the AI/ML community [18, 38, 42]. An increasing presence of publications at major conferences/journals concerned with the integration of causality with AI/ML (including [24, 27–29, 32, 54] to mention a select few) suggests a growing subfield that sets a consensus on *causal* AI/ML as promising paradigm for next-generation systems. Still, as the difficulty of the integration with otherwise prominent success stories of deep learning, such as computer vision, becomes apparent, countering opinions start to speak out against causal AI/ML [7]. In this work, we take the former perspective *pro* causal AI/ML. We argue that the questions around causality can fuel research also on questions of recent debates such as how much 'real' progress towards AGI has been made since the advent of large scale models such as BERT [16], GPT-3 [12], DALL-E [40].

The following block paragraph serves as a summary of an example of such a recent debate on whether 'just scaling' these models is a sufficient condition for progression towards AGI:

With the rise of large scale models such as BERT, GPT-3, DALL-E, AI history suggests to repeat itself as arguably impressive text generation and image synthesis results foster different opinions in the community in terms of interpretation regarding the progression of the field as a whole towards the grand goal of AGI (key references involve [30, 31] that sparked intense discussions amongst notable researchers via social networks; also for reference, a short treatise that discussed patterns in the history of AI research observes: "early, dramatic success followed

^{*}Both authors contributed equally to this research.

Authors' addresses: Moritz Willig, moritz.willig@cs.tu-darmstadt.de; Matej Zečević, matej.zecevic@tu-darmstadt.de, Technical University of Darmstadt, Darmstadt, Hessen, Germany; Devendra Singh Dhami, devendra.dhami@tu-darmstadt.de, Technical University of Darmstadt and hessian.AI, Darmstadt, Hessen, Germany; Kristian Kersting, Technical University of Darmstadt, hessian.AI and DFKI, Darmstadt, Hessen, Germany.

^{2023.} Manuscript submitted to ACM

by sudden unexpected difficulties." [14]). Some even speak of *foundation* models [8] to account for the "emerging paradigm" of models that provide a base from which task-specific models are derived through adaptation. The emergence of what seem to be the two dominant clusters of opinions within the discussion around foundation models is characterized by researchers who recognize said models as significant progression towards AGI and those who do *not*. For the former group of the observed results act as corroborating evidence for the *scaling hypothesis* [11, 46] which captures that the idea of emergent properties as a result of scaling neural network in terms of parameters and data, thereby, rooting parts of the overarching idea in results from neuroscience that suggest the human brain to 'just' be a scaled up primate brain [21]. The other popular opinion is that the achieved results act as a mere reflection of the sheer scale of data and parameters, put differently "the methods are old" and their lack of interpretability and reasoning capabilities will remain persistent. While many researchers voiced there opinions and thoughts, one notable comment came from Judea Pearl who announced his alliance with the latter position via social media, stating "These models are castles in the air. They have no foundations whatsoever." discrediting the models for lacking any identifiable notion to causality.

It is clear how resolving these discussion through scientific inquiry is crucial for the AI/ML community. The question of the whether or not Pearl's statement is true was one of the seeds that grew in to the present paper.

Therefore, in this paper, we investigate whether current foundation models are such "castles in the air."

We identify the key problem of ongoing debates to lie in the scale of data and parameters that only further cement the inherently black-box nature of the base models. Therefore, to answer whether such "foundation models" have made progress towards AGI and to give reason onto why causal AI/ML could be a milestone, it seems to suffice to ask and investigate the question of the extent to which *foundation models can talk causality*. For the sake of simplicity, we will take the "talking" literally and focus on LLMs in this work while leaving general foundation models (e.g. image-based models) for future work.¹

The paper is structured as follows: we investigate how 'causal' LLMs are by first formalizing our key hypothesis on "correlations of causal facts" using Pearl's language, introducing necessary new ideas along the way with illustrative examples, posing our key theoretical contribution. Following that, we provide an empirical analysis on the causal prowess of current LLMs and discuss the results in lights of our theoretical groundwork from the beginning, posing our second contribution.

For reproduction purposes we make our code repository for the empirical part publicly available.²

2 INFORMAL SUMMARY OF THE MAIN IDEA OF THE PAPER

LLMs are transformer-based models [51] and it is clear how they are not parameterized variants of SCMs such as the neural ones presented in [53] as they do not specify structural equations mimicking causal mechanisms. Nonetheless, they do seem to answer causal questions right sometimes. Our explanation for this is that they are not only 'stochastic parrots' as already suggested by Bender et al. [6] but sometimes also 'causal parrots' since they will also encounter correlations over causal facts during training in their vast oceans of textual data. Essentially, you can train LLMs to tell you anything that sounds nice to you, in this case 'nice' meaning right causal answers, but that is all. We illustrate this idea in Fig.1 where we show the example of the physical concept that we know as altitude and how it is being causal of the concept of temperature-in the sense that there is a physical mechanism that leads to a temperature

¹However, for observations that we believe should also hold more generally for foundation models (like the concept of "correlations of causal facts" that we will introduce in this work) we will use the term foundation model instead of writing LLM.

¹⁰³ ²https://anonymous.4open.science/r/causalParrots-75D9/

¹⁰⁴ Manuscript submitted to ACM





Fig. 1. **Same Implication, Different Representations.** When we consider the causal relationship between altitude (A) and temperature (T), then it is apparent that given the laws of physics we have an increase in altitude leading to a decrease in temperature. Graphically we can depict the relationship as $A \rightarrow T$, whereas the actual 'increase-decrease' relationship can only be specified through the SCM formalism with its structural equations, that is some f such that T = f(A, U) where U are exogenous variables. The ground truth SCM underlying our laws of physics generates observational data in the form of numerical tuples (a, t) as seen on the left scatter plot. To infer the casual relation, we can resort to algorithms for causal discovery. However, crucially, the same knowledge achieved through such induction can be represented within *text* for 'free' as one simply recites the Wikipedia article found on the right. While the article on the right is correct, and thus represents a fact about the actual world, there is no such guarantee for arbitrary other texts. That is, a model that simply obtains its knowledge from various Wikipedia statements will also learn untrue statements, statements that are not facts, thus explaining behavior that is correct sometimes and wrong other times.

increase/decrease with a corresponding change in altitude-is a fact of causality in our physical reality that can be 127 represented and therefore also learned in different ways. One way, the way that we usually consider in ML, is through 128 induction on physical measurements. We have actual data points of different altitude-temperature pairs (maybe recorded 129 from different cities and mountain ranges) and infer the underlying relationship. For this we could fit a linear causal 130 131 model with non-gaussian noise (LiNGAM; see the original work for reference [44]) and see that it explains our data 132 perfectly, concluding that altitude causes temperature. However, this conclusion is arguably nothing new, as most 133 people would agree, and this is partly so because such obtained knowledge has been embedded as textual articles into 134 encyclopedias such as Wikipedia, which are freely accessible. Turns out that not just humans but also LLMs learn 135 136 from precisely such textual data that already contain the results from the original physics experiments with 'real' data. 137 For instance, 'The Pile' data set which was used for training OPT [56] compromised over 6 GiB of textual data from 138 Wikipedia [17]. Both physical measurements and the facts we find on Wikipedia are a consequence of the causal reality 139 that altitude causally influences temperature and reversely they both imply this same causal model, however, they are 140 141 fundamentally different forms of representation in that learning from the physical measurements can be argued to 142 be what we mean by 'understanding something', whereas simply reading up on the textual article lacks exactly that 143 component. Indeed, this discussion then turns philosophical onto which we will add a few more comments. 144

While the philosophical argument introduced by Searle [43] is essentially the same as the stochastic/causal parrot 145 146 argument minus the probabilistic view, another compelling comparison of LLMs is that of the Plato's allegory of the 147 cave Plato [39] in which the question is raised to which extent one can learn about the real world's functioning by just 148 observing the shadows of its objects. For this allegory, we could see the LLM act as the cave where one can observe 149 some causality in the form of correct answers (the 'shadows') but the question is raised whether this would in fact be 150 151 actual causality (the 'real world'). Even if we consider LLMs to be universal in the same sense that we have proven that 152 neural networks are universal function approximators Cybenko [15] that does not imply that it is easy to make them 153 causal and therefore even if they can make use of correlations exposed by meta SCM talking about causal facts, then 154 still the LLM would be required to be exposed to an infinite amount of data from such universal, meta SCM. 155

Manuscript submitted to ACM

157 158

172

173

174

175 176

177

178 179

180

181

182

183 184

185 186

187

188

189

190 191

192

193 194

195

196

197

198

3 FORMALIZING "CORRELATIONS OF CAUSAL FACTS"

"Correlation does not imply causation," goes the famous saying (see Aldrich [2], Pearl [36]), that accounts for the fact 159 that following Reichenbach's common cause principle a correlation between two variables might be either because 160 161 one is causing the other, or because there is a third variable causing both [41]. To infer the 'actual' causation³ within 162 the system of interest, we might resort to manipulating the system, as another fomous saying suggests "No causation 163 without manipulation" [23]. A celebrated victory of Pearl's notion to causality is the causal hiearchy theorem (CHT) 164 which guarantees that purely observational data collected from a system can not be used to uniquely determine causal 165 166 statements, when no other causal assumptions are available [4]. The CHT certainly seems to imply that no matter 167 how much we scale our foundation models (in terms of data and parameters), we will never be able to perform causal 168 inference. In a nutshell, the CHT seems to disprove the scaling hypothesis. Or does it? In this work, we argue that 169 foundation models might be exploiting a "loop hole" in the CHT⁴. Namely, what happens if the causal assumptions 170 (which are required, by the CHT, for causal inference) are represented in observational data itself? 171

We start by providing the definition of a Structural Causal Model (SCM). We will follow the broadest notion of SCM that still enjoys properties such as unique solvability, marginalization and an intuitive graphical underpinning. The class is that of *simple* SCM as first introduced by Bongers et al. [10]. They extend typical acyclic, semi-Markovian⁵ SCM with cycles that offer a unique solution⁶. Formally, we have:

DEFINITION 1. A simple Structural Causal Model (SCM) is a tuple $\mathcal{M} := (I, \mathcal{J}, X, \mathcal{E}, f, \mathbb{P}_{\mathcal{E}})$ where I, \mathcal{J} are disjoint, finite index sets⁷ for endo- and exogenous variables respectively, $X = \prod_{i \in I} X_i, \mathcal{E} = \prod_{j \in \mathcal{J}} \mathcal{E}_i$ are products of the domains of the variables where each domain is a standard measurable space⁸, f are the structural equations for which X = f(X, E)for random variables X, E, and finally $\mathbb{P}_{\mathcal{E}}$ denotes the exogenous distribution. \mathcal{M} is further uniquely solvable⁹ for every subset $O \subseteq I$. The graph implied by \mathcal{M} is denoted as $\mathcal{G}(\mathcal{M})$.

The graph of a simple SCM is drawn by considering the parent-child relationships implied by the structural equations. We call a variable *k* parent of some variable *i* if there is exists no measurable function that can match the actual structural equation that computes the values of *i* while not using *k* as argument. For simple SCM we then end up having three basic types of *causal* relationships that can exist for any variable pair (i, j) of $\mathcal{G}(\mathcal{M})$: (1) if there is a direct edge $i \rightarrow j \in \mathcal{G}(\mathcal{M})$, then *i* is called a direct cause of *j*, (2) if there is a directed path $i \rightarrow \cdots \rightarrow j \in \mathcal{G}(\mathcal{M})$, then *i* is simply a cause, and (3) if there is a bidirected edge $i \leftrightarrow j \in \mathcal{G}(\mathcal{M})$, then *i* and *j* are confounded. In the following, we will provide an example of a simple SCM based on the 'classical' setting described in Fig.1 using the above definition.

EXAMPLE 1 ('CLASSICAL SETTING'). Let X, Y, Z be the random variables whose values describe "Countries's Annual Per Capita Chocolate Consumption", "Number of Nobel Laureates per 10 Million Population" and the "Gross Domestic Product" for any given country respectively. The data observed in Messerli [34] suggested a significant linear correlation (r = 0.791, p < 0.0001) for (X, Y) with Switzerland being the top-performer (X > 10, Y > 30). Note that Z was not observed

 $[\]frac{199}{3}$ For a rigorous treatment of different notions of what consitutes an "actual cause" consider the seminal work of Halpern [20].

²⁰⁰ ⁴Or rather, it is a *subtle* detail that might easily be forgotten.

²⁰¹ ⁵SCMs that allow for latent confounding, that is, the exogenous terms need not be independent.

⁶However, these do unfortunately exclude self cycles and many other circular relationships.

⁷For examples later on, we might simply equate indices to letters when considered more suitable for presentation, that is, instead of having $I := 3 = \{1, 2, 3\}$ we write $I := \{X, Y, Z\}$.

⁸See Def.F.1 of [10] that gives a rigorous, albeit very technical definition. Typically, definitions on SCMs found in the literature are rather hand-wavy about what the 'variables' are. For instance Bareinboim et al. [5] (Def.1) simply refers to 'variables' in the most general sense, while Pearl [36] is more

specific in talking about 'random variables'. However, probability theorists often restrict themselves to separable completely metrizable spaces also

known as Polish spaces, which avoid pathological situations in which the considered spaces' sizes are uncountable, that is, greater/equal 2^{N_0} .

²⁰⁷ ⁹See Def.3.1. in [10].

²⁰⁸ Manuscript submitted to ACM

217 218

219

220

221

226

227

228 229

230

231 232

233

234

235

236 237

238 239

240

241

242 243

244 245

246

247 248

260

but clearly serves as a reasonable explanation for the underlying causal system. As postulated by Reichenbach's common 209 210 cause principle [41], if the correlation does not imply a direct causation, then there must exist a common cause, which we'd 211 choose as Z for this particular example. While we do not know the true SCM accountable for X, Y, Z, we can approximate 212 the data from Messerli [34] reasonably well using the following simple SCM $\mathcal{M}_1 := (\{X, Y, Z\}, 3, \mathbb{R}^3, \mathbb{R}^3, f, \mathbb{P}_{\mathbb{R}^3})$ with f213 214 being defined as

$$\mathbf{f} := \{ X = f_1(E_1) = E_1, \quad Y = f_2(X, E_2) = 2 \cdot X + E_2, \quad Z = f_3(E_3) = E_3 \}^{10}.$$
(1)

The structural equations in (1) make \mathcal{M}_1 a linear additive noise model that approximate [34] reasonably well. It is clear how the observed correlation in this case corresponds to a direct causation according to f since X is a parent of Y and Z simply an independent variable. However, we are certain that chocolate consumption does not cause more Nobel laureates anywhere, so M_1 is not reasonable in predicting our real world expectations. Following Reichenbach's postulate, it would be more reasonable to use an alternate SCM $\mathcal{M}_2 := (\{X, Y\}, \mathbf{3}, \mathbb{R}^2, \mathbb{R}^3, \mathbf{f'}, \mathbb{P}_{\mathbb{R}^3})$ with $\mathbf{f'}$ being defined as

$$f' := \{ X = f'_1(E_1, E_3) = E_3 + E_1, \quad Y = f'_2(E_2, E_3) = 2 \cdot E_3 + E_2 \}.$$
(2)

This second SCM M_2 would now correspond better to both (1) the actual data observed in [34] since Z was never observed and is modelled as exogenous variable E_3 implicitly and (2) to our real-world intuition since there is a bidirected edge $X \leftrightarrow Y \in \mathcal{G}(\mathcal{M}_2)$ with E_3 being the underlying confounder.

Example 1 serves to show how the rather abstract definition of an SCM can be made tangible to communicate what we believe about our observed data and more so the underlying data generating process. Previously, we have defined a (direct) cause as directed (1-)path in the causal graph, however, we have not discussed why we call such an edge in the implied graph of an SCM a 'cause'. The reasoning behind the naming turns out is an important insight that we will use soon to develop our idea of "correlation of causal facts." But first, we need to briefly talk about Pearl's Causal Hierarchy (PCH) which defines three (symbolic) languages $\mathcal{L}_1, \mathcal{L}_2, \mathcal{L}_3$ with increasingly expressive quantities. For this, we make use of the definitions proposed by Bareinboim et al. [5].

DEFINITION 2. The Pearl's Causal Hierarchy (PCH) consists of three (symbolic) languages $\mathcal{L}_1, \mathcal{L}_2, \mathcal{L}_3$. The terms are called: (1) observational $P(\mathbf{Y} = \mathbf{y}) \in \mathcal{L}_1^{11}$, (2) interventional $P(\mathbf{Y}_{\mathbf{x}} = \mathbf{y}) \in \mathcal{L}_2$, and (3) counterfactual $P(\mathbf{Y}_{\mathbf{x}} = \mathbf{y}, \dots, \mathbf{Z}_{\mathbf{w}} = \mathbf{y})$ $z) \in \mathcal{L}_3$. Where $P(Y_x = y)$ denotes the probability of Y being y were X to have values x^{12} .

Now we can return to clarifying why we call certain paths in the causal graph 'causal'. We state our insight as:

INSIGHT 1. Let \mathcal{M} be some SCM. Knowledge about the structural equations and the causal graph of \mathcal{M} is knowledge about answering \mathcal{L}_3 and \mathcal{L}_2 queries in \mathcal{M} respectively.

Returning to Example 1, it is clear how knowing the actual parameterizations of the functions f'_1, f'_2, f'_3 allows us to 249 answer \mathcal{L}_3 queries and not knowing the actual parameterizations but at least knowing that for instance variable Z is a 250 251 parent of X and therefore a necessary argument in f'_1 is sufficient for answering \mathcal{L}_2 queries. Put differently, we can 252 call a direct edge $Z \to X \in \mathcal{G}(\mathcal{M}_2)$ 'causal' since it is a \mathcal{L}_2 fact¹³. We can rephrase this statement differently to make 253

¹⁰Note how these equations make the assumption that all the variables are measured using the *same scale*. This is of course not a necessary assumption, 254 however, since this example serves only illustrative purposes for clarifying the ideas around SCMs, and we do not train an actual SCM on the data from 255 [34], the chosen presentation is reasonable. 256

¹¹We use $P(\mathbf{X} = \mathbf{x})$ to simply denote the probability value of random variable **X** obtaining value **x**.

¹²Not to be confused with regular conditionals $P(\mathbf{Y} = \mathbf{y} \mid \mathbf{X} = \mathbf{x}) \in \mathcal{L}_1$ that are read as "the probability of Y being y were X observed as X". 257

 $^{^{13}}$ How we actually acquire such \mathcal{L}_2 (or even \mathcal{L}_3) facts is a different question altogether. It stands at the heart of the discipline of *causal discovery*. We call 258 \mathcal{L}_2 'interventional' since we can use experimentation in the form of 'surgical intervention' on admissible variables measured in the given task to discern 259 causation from mere correlation.

clear where our overall idea is leading to. Namely, some machine learning model is 'causal' w.r.t. some query if it can answer that query with the right \mathcal{L}_2 fact. More intriguingly, it does not matter where that \mathcal{L}_2 fact comes from since the formulation is independent of whether the model learns the fact and simply requires that the model knows about the fact. We state our second key insight as:

INSIGHT 2. The 'variables' of SCMs are not restricted to 'natural' concepts such as "Chocolate consumption" or "Number of Nobel laureates" (see Ex.1), they can be 'meta' concepts involving causal facts, that is, knowledge about \mathcal{L}_2 and \mathcal{L}_3 .

For our argument the existence of data describing \mathcal{L}_2 facts is already sufficient, therefore, for ease of presentation we will focus on those specifically in the following although we see no difficulties (as of writing) in extending to \mathcal{L}_3 . We

call these concepts 'meta' since they are one level above 'regular', simple SCM in the sense that they encode information about answering causal questions in another SCM. To make this idea more formal, we define 'meta' SCM as follows:

DEFINITION 3. Let \mathcal{M}_1 and \mathcal{M}_2 be two SCMs such that the observational distribution of \mathcal{M}_2 denoted $\mathcal{L}_1(\mathcal{M}_2)$ can answer queries w.r.t. the interventional distributions of \mathcal{M}_1 denoted $\mathcal{L}_2(\mathcal{M}_1)$, then \mathcal{M}_2 is called meta (w.r.t. \mathcal{M}_1).

In the following, we construct an example of such a meta SCM analogue to the 'classical setting' from Example 1.

EXAMPLE 2 ('META SETTING'). As before, let X, Y, Z denote the natural concept random variables from Example 1 for the chocolate-Nobel data from Messerli [34] and set $\mathcal{M}_1 := (\{X, Y\}, \mathbf{3}, \{0, 1\}^2, \{0, 1\}^3, \mathbf{f}, \mathbb{P}_{\{0, 1\}^3})$ with $\mathbf{f} := \{X = f_1(E_1, E_3) = f_1(E_1, E_3) = f_2(E_1, E_3)$ $E_3 \wedge E_1, Y = f_2(E_2, E_3) = E_3 \wedge E_2$ where X = 1, Y = 1 denote 'high' chocolate consumption and number of Nobel laureates respectively, and vice versa for X = 0, $Y = 0^{14}$. In this example, we intend on answering an arbitrary example query from $\mathcal{L}_2(\mathcal{M}_1)$ like for instance $P(Y_{X \leftarrow 1} = 1)$, that is, the probability of a high number of Nobel laureates if the given chocolate consumption were to be high. Clearly, since $X \leftrightarrow Y \in \mathcal{G}(\mathcal{M}_1)$, we expect $P(Y_{X \leftarrow 1} = 1) = P(Y = 1)$ since intervening on X will not change Y. For the sake of the example let's assume fair conflips as exogenous distributions, $\mathbb{P}_{\mathcal{E}_i} := \mathcal{B}(1/2)$, then $P(Y_{X \leftarrow 1} = 1) = 1/4$. Next, we need to show the existence of some SCM M_2 that is meta to M_1 , that is, which will answer $P(Y_{X \leftarrow 1} = 1)$ using $\mathcal{L}_1(\mathcal{M}_2)$. This is easy, as we can define

$$\mathcal{M}_{2} := (\{W\}, \emptyset, \{0, 1\}^{3 \times 3}, \emptyset, f', \mathbb{P}_{\emptyset})$$
(3)

where $W := \mathcal{G}(\mathcal{M}_1)$ is a random variable that describes a causal graph over X, Y, Z as an adjacency matrix¹⁵ and therefore

$$\boldsymbol{f'} := \{ W = f_1'(W) = \begin{pmatrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 1 & 0 \end{pmatrix} \}.$$
(4)

W encodes the causal graph $X \leftarrow Z \rightarrow Y$. It is now a simple exercise to show that M_2 is indeed meta. We want to find the a in $P(Y_{X \leftarrow 1} = 1) = a$. We know that $P(Y = 1) = \frac{1}{4}$ and using $\mathcal{L}_1(\mathcal{M}_2) = W$ we further know that X is **not** a parent of Y since $X \to Y \notin W = \mathcal{G}(\mathcal{M}_1)$. Therefore, we conclude that $P(Y_{X \leftarrow 1} = 1) = P(Y = 1) = 1/4 = a$.

302 At this point, it is worthwhile noting that while every SCM is in fact a causal model it is not the case that every 303 SCM is a reasonable causal model for our physical world or rather for what is being modelled. While the meta SCM in the above case can also be viewed even as a simple SCM which retains even its causal interpretation on graphs, 305 306 it is surely not a 'regular' SCM particularly because it is in direct reference of another SCM. We have achieved two things in the above example (1) shown the existence of meta SCM and (2) shown how they can actually answer concrete 308 causal facts about another SCM. While M_2 from the above example may seem artificial in the sense that it is tailored to 309

6

261 262

263

264

265 266

267

268 269

270 271

272

273

274 275

276

277 278

279 280

281

282 283

284

285

286

287 288

289

290

295 296

297

298 299

300

301

304

³¹⁰ ¹⁴We adapted the SCM to be binary for ease of calculation in the example, the definition holds w.l.o.g. for any simple SCM.

³¹¹ ¹⁵A cell $\hat{A}_{ij} = 1$ indicates the direct edge $i \rightarrow j$.

³¹² Manuscript submitted to ACM

315

316

317 318

319

320 321

322

337

338 339

340

341

342 343

344

345

346

347 348

349 350

351

364

clarity for the example.

 \mathcal{M}_1 , the example acts as intended since we could clearly define a suitable meta SCM for our illustration's purposes.¹⁶ After having proven the existence of meta SCM that can encode \mathcal{L}_2 facts (of another SCM) within their observational distribution, hence making them meta, we are finally ready to provide our main running hypothesis for this work that we believe to hold, namely correlations of causal facts.

CONJECTURE 3.1 (CORRELATION OF CAUSAL FACTS (CCF)). As before, let M_1 be some SCM and M_2 a respective meta SCM. Further let $Q \in \mathcal{L}_2(\mathcal{M}_1)$ and $A \in \mathcal{L}_1(\mathcal{M}_2)$ be causal queries with their respective answers and f denotes the LLM's predictive model. Then we have: $f(Q) = A \iff f(Q)$ minimizes training error.

323 In words, Conj.3.1 suggests that in all the cases where the LLM does provide the right causal answer to a causal query, 324 then it is only because (i) this fact was observed in the training data and (ii) the correlation with the query is optimal 325 from the perspective of the training objective. This formulation also leaves room for the fact that LLMs being right 326 will not always (and maybe not even most of the times) be the case. Furthermore, it provides another very intriguing 327 328 observation in that an LLM, although we conjecture it to essentially be involved with at least two SCM (the one 329 modelling the phenomenon of interest and the meta one modelling the first model) during training, it is not causal and 330 cannot be (up to answering some subset of causal questions right). To prove our conjecture, we would require three 331 things: (i) identify meta SCM that generated the causal facts in the LLMs training data, (ii) show that the causal facts are 332 333 the best (in terms of training loss) answers to the causal queries and (iii) show that LLMs indeed give the causal facts 334 as answers to the causal question. Since (i-ii) are arguably infeasible or at least fiendishly difficult to validate, in the 335 following, we will focus our initial empirical analysis on (iii) and see judge the causal inference capabilities of LLMs. 336

4 TESTING FOR CAUSAL KNOWLEDGE IN LARGE LANGUAGE MODELS

We evaluate three publicly accessible LLMs: OpenAI's GPT-3 [12], AlephAlpha's Luminous [3], and Meta's OPT [56]. All models are transformer based architectures [51] trained at scale, qualifying them as LLMs (see Bommasani et al. [8]). Our analysis primarily investigates three different questions that belong to part (iii) of CCF as described earlier, namely: How do LLMs perform ..

- .. in "common sense" settings like reasoning or intuitive physics?
- .. in settings where the causal graph is (partially) known?
- ..when using their embeddings of knowledge base facts?

In the following, we conduct various experiments to answer these questions.

4.1 Methods and Results for "Common Sense" Inference Tasks

We argue that "common sense" reasoning tasks either involving some basic propositional logic or intuitive physics (as 352 353 reference consider for instance Tenenbaum et al. [48]) are reasonable settings in which we can expect CCF to hold. 354 For propositional logic we consider 20 different questions such as for example "If A causes B and B causes C does A 355 cause C?". These questions are simply fed as prompt to the respective LLM. In this setup no prior world knowledge is 356 required, other than being able to link together propositions. In the simplest propositional chain we provide the model 357 358 with a chain of three variables, linked by two propositions: "A causes B", "B causes C." We ask the model "Does A cause 359 C?". This setup provides a fairly simple and easy to control experiment. We vary the length of the chain up to n = 10360 variables: X_1 causes X_2 ... and X_{n-1} causes X_n . To prevent the model from simply pattern matching the first and last 361 362 16 We could have chosen W to represent the space of all directed acyclic graphs (DAG) and eventually found the same conclusion, albeit at the cost of 363

Manuscript submitted to ACM

365								Int	uitive	Phy	sics				
366		Rolling	(8)	Suj	ppor	t (8)	Co	llisic	ons (4) S	eesaw (4)	Wei	ghts (5)	Tools (7)	Accuracy
367	GPT-3	6		5			4			2		2		3	61.11%
368	Luminous	1		0			0			1		1		2	11.11%
369	OPT	2		0			1			0		0		4	19.44%
370		C	ausa	l Cha	ains	(Basi	c Pro	p. L	ogic)						
371		N=2	3	4	5	6	7	8	9	10	Subchain	s (4)	Randor	nized (7)	Accuracy
372	GPT-3		1	1	1			1		1	2		2		45.00%
373	Luminous	1				1	1	1	1		1		4		50.00%
374	OPT		1			1					0		2		20.00%

Table 1. (top) Querying LLMs for questions about intuitive physics. GPT-3 performs reasonably well across all queried areas, while
 OPT performs better on the Tools subtask. (bottom) Querying LLMs for propositional logic on causal chains. GPT-3 and Luminous
 perform equally on causal chains of varying length. Luminous outperforms the other models when asked about causal chains with
 randomized chain ordering or variable names. Overall LLMs display a mixed performance.

variable we ask for sub-chains (e.g. "Does *B* cause *E*?") that leave out the initial or last part of the chain. Additionally, we
 might also switch the order of propositions or exchange the alphabetical ordered variable names with random letters.

Results for this experiment are shown in Table 1. A complete list of questions and answers can be found in the 383 appendix. We observe that GPT can handle shorter chains up to and fails starts to fail at n = 6 variables. Contrary to 384 this, Luminous start to answer correctly at n = 6 variables but fails at shorter chains. OPT only answers two queries 385 386 correctly. For most of our evaluations we observe the tendency of OPT to repeat the given query text and not answering 387 the prompt. For all three models we see a decline in performance once we leave the standard $X_1...X_n$ setting and 388 start to query for sub-chains and randomize proposition order. Only Luminous keeps reasonable performance for the 389 'Randomized' queries. From our observations we conclude that current models still need special fine tuning such as in 390 391 [55] to be applicable to inference tasks with given information.

392 For intuitive physics we consider 36 questions such as for example "A ball is placed on a table and rolls off. What 393 does this tell us about the table?". As we are not able to feed images depicting the physical scenes into the LLMs, 394 395 we are resorting to textual descriptions. We provide idealized descriptions of scenes that query the understanding 396 about physical mechanisms, such as balls rolling off tiled surfaces or objects bouncing away in a collision. Using 397 textual descriptions leads to ambiguity as weights, sizes and relative positions of the objects might not be described. To 398 compensate for this uncertainty, we choose simple setups that might even occur in text books and count all *plausible* 399 outcomes given by the LLMs (as manually evaluated by us) as correct. 400

401 Again, results for this experiment are shown in Table 1. The most striking observation that caught us were some 402 parallels to human nature of reasoning. In the 'Weights' cluster of questions we asked the LLMs to answer variations of 403 the well-known trick question "What is heavier: A kilogram of metal or a kilogram of feathers?". For this questions all 404 models wrongly opt for the material with higher specific weight and answer "A kilogram of metal is heavier than a 405 406 kilogram of feathers." However, when asked "['A kilogram of metal is heavier than a kilogram of feathers'] is what 407 most people say, but in reality", GPT correctly answers "They weigh the same." arguably the same way many humans 408 intuitively do! This is interesting from a factual point of view, as the extended formulation contains no new factual 409 information. A pure propositional reasoner would still have come to the same answer independent of the truth value of 410 411 the question, as the additional information about what "most people say" does not factually contribute to comparing 412 two different weights. It is rather a hint on the meta level that is often used in human communication to indicate 413 the presence of a pitfall or trick question, and therefore not to choose the 'obvious' answer. In the presence of such a 414 hint GPT does not just negate its previous given answer but arrives at the fact that both weigh the same. Seeing LLM 415 416 Manuscript submitted to ACM

considering these meta level hints might help with future algorithms that integrate LLMs into their processes. Providing
 textual hints might be comparable to switching between different heuristics as it is done in today's classical algorithms.

4.2 Methods and Results for Causal Discovery on Ground Truth

420

421

422 Previously, we have seen that LLMs do not necessarily perform well on tasks where they are not just required to 423 know a fact but also to reason with that fact. For the following experiment, the LLM will not be required to reason 424 with facts but it will be sufficient to simply recall the right facts. Therefore, we can now simply resort to causal 425 discovery benchmarks for evaluating whether LLMs can recover causal graphs when queried accordingly. We propose 426 427 a naïve causal discovery approach for LLMs to this end. Fig.2 provides a schematic overview of the naïve structure 428 discovery procedure that we propose. To account for stability and reproducibility, we present different wordings to the 429 queries (synonymous formulations) and disable parameters that induce randomness (e.g. temperature), respectively. 430

It is important to note that the proposed naïve structure discovery 431 432 procedure is not a proper induction method in the classical sense 433 as it does not use actual data as input to perform the inferences (all 434 the possible inferences are established upon training completion). In 435 that sense, LLMs behave much like humans, who simply recall that 436 437 "a higher altitude means a lower temperature" than to look at actual 438 data recordings of altitude and temperature (and other variables) to 439 perform the causal inference. As anticipated, the LLM thereby also 440 inherits natural language ambiguities. To give an example, even if 441 the LLM is prompted with an additional "Answer with Yes or No" 442 443 the LLM is not constrained to oblige. We use five different query 444 wordings (or formulations) such as "Are X and Y causally related?" 445



Fig. 2. Naïve Causal Discovery with LLMs.

or "Does X cause Y?" (see appendix for full list). The first three of which are classified as symmetric queries since the 446 447 expected answer is a mere association X-Y and the last two wordings classify as *asymmetric* accordingly i.e., we expect 448 either $X \to Y$ or $X \leftarrow Y$ (in the case of an existing relation). For any answer given by the LLMs we automatically classify 449 answers starting with 'Yes' or 'No' accordingly and manually label the remaining ones. For any data set containing 450 N variables there are $\binom{N}{2}$ possible edges. Since the edges are directed we get twice the amount of queries, one query 451 for each direction, multiplied by the number of query wordings Q = 5 and arrive at $2 * \binom{N}{2} * Q$ queries per data set. 452 453 We consider publicly available data sets that propose a 'ground truth' causal graph (which depicts the data generating 454 process). We consider six data sets: altitude (A; Mooij et al. [35]), health (H; Zečević et al. [54]), recovery (R; Charig 455 et al. [13]), driving (D; synthetic), cancer (C) and earthquake (E) both [26]. For all data sets we query the LLM with all 456 possible combinations of edges between any two variables. For our data sets, we get 10 questions for Altitude, 100 for 457 458 Cancer, 60 for Health, 30 for Driving, 100 for Earthquake and 30 for Recovery respectively. 459

After querying all edges for all of data sets we then compare each of the obtained graphs to the respective ground 460 truth using different metric. In total we discuss two key observations. Table 2 presents the results of our experiment. 461 462 We present metrics to measure different aspects of the LLM predictions. We measure fitness to the causal graph using 463 Structural Intervention Distance (SID; [37]) and Structural Hamming Distance (SHD; [1, 49]). Machine Learning (ML) 464 applications might care more about capturing all relevant connections and less about including too many. To reflect 465 this we report the F_1 score. Furthermore, we inspect individual statistics on the *edges* of the predicted graph when 466 467 wording of the queries changes. We define sparsity as the number of predicted edges in relation to the maximum 468 Manuscript submitted to ACM

	Metric	Altitude	Health	Driving	Recovery	Cancer	Earthquake	LLM
.1		$0.80_{\pm 0.40}$	$7.20_{\pm 0.75}$	$3.00_{\pm 0.89}$	$4.00_{\pm 1.79}$	$11.80_{\pm 4.66}$	$11.40_{\pm 1.50}$	GPT-3
ld b	$SID \downarrow$	$1.20_{\pm 0.98}$	$10.60_{\pm 1.85}$	$6.00_{\pm 0.00}$	$5.40_{\pm 1.20}$	$11.40_{\pm 3.07}$	$16.00_{\pm 3.63}$	Luminous
G		$1.60_{\pm 0.80}$	$10.80_{\pm 2.40}$	$5.00_{\pm 1.26}$	$5.80_{\pm 0.40}$	$16.80_{\pm 1.94}$	$15.60_{\pm 5.95}$	OPT
sal		$0.80_{\pm 0.40}$	$4.00{\scriptstyle\pm0.63}$	$2.60_{\pm 0.49}$	$2.20{\scriptstyle \pm 0.40}$	$7.00_{\pm 1.41}$	$4.60{\scriptstyle \pm 0.80}$	GPT-3
Cau	SHD \downarrow	$0.60_{\pm 0.49}$	$7.00_{\pm 1.10}$	$4.20_{\pm 0.40}$	$3.40_{\pm 0.80}$	$10.00_{\pm 3.52}$	$5.60_{\pm 1.62}$	Luminous
0		$0.80_{\pm 0.40}$	$7.40_{\pm 1.20}$	$3.40_{\pm 1.20}$	$4.00_{\pm 0.00}$	$13.20_{\pm 1.60}$	$8.60_{\pm 3.01}$	OPT
	F_1 Score \uparrow	$0.20_{\pm 0.40}$	$0.47_{\pm 0.14}$	$0.11_{\pm 0.23}$	$0.27_{\pm 0.33}$	$0.35_{\pm 0.11}$	$0.12_{\pm 0.15}$	GPT-3
WI		$0.80_{\pm 0.16}$	$0.41_{\pm 0.21}$	$0.46_{\pm 0.09}$	$0.55{\scriptstyle \pm 0.07}$	$0.40_{\pm 0.13}$	$0.40_{\pm 0.04}$	Luminous
		$0.73_{\pm 0.13}$	$0.52_{\pm 0.05}$	$\textbf{0.53}_{\pm 0.15}$	$0.47_{\pm 0.07}$	$0.35_{\pm0.03}$	$0.47_{\pm 0.07}$	OPT
		$0.90_{\pm 0.20}$	$0.63_{\pm 0.28}$	$0.77_{\pm 0.31}$	$0.70_{\pm 0.31}$	$0.65_{\pm 0.16}$	$0.93_{\pm 0.07}$	GPT-3
	Sparsity	$0.20_{\pm 0.24}$	$0.22_{\pm 0.35}$	$0.03_{\pm 0.07}$	$0.10_{\pm 0.13}$	$0.40_{\pm 0.16}$	$0.74_{\pm 0.12}$	Luminous
ges		$0.10_{\pm 0.20}$	$0.05_{\pm 0.10}$	$0.17_{\pm 0.21}$	$0.07_{\pm 0.13}$	$0.18_{\pm 0.12}$	$0.41_{\pm 0.18}$	OPT
Ed_{j}		0.50	0.62	0.33	0.50	0.69	0.00	GPT-3
	ADS ↑	1.00	0.53	0.17	0.17	0.38	0.26	Luminous
		0.50	0.25	0.25	0.33	0.28	0.47	OPT

Table 2. Comparing LLMs prediction to existing ground truth causal structures. The metrics concerned with the causal graph structure (SID, SHD) reveal a closer match of GPT-3 predictions to the ground truth causal structures than for the other LLMs. High F_1 Scores and low sparsity indicate densely connected graph prediction by Luminous and OPT. This can be desired for ML applications. The ADS reveals that all LLMs increase their decisiveness on edge directions when querying with asymmetric sentence templates.

number of possible edges.¹⁷ Intuitively, this metric measures the percentage of maximally predictable edges and gives an insight on whether a LLM tends to predict sparse (Sparsity ≈ 1) or densely connected graphs (Sparsity ≈ 0). Additionally we want to measure the amount of directed edges within a graph. While SCM are not restricted to DAG structures, most classical causal data sets are modeled using DAGs. We further introduce a novel metric based on the model's decisiveness d which is defined as the percentage of directed edges within the total number of predicted edges. The exact procedure is shown in Algorithm 1 in the appendix. When comparing two graph structures we define Δd to be the change in decisiveness between two predictions: $\Delta d(A, B) := d(B) - d(A)$. A positive Δd value indicates an increase in directed edges when switching from A to B and vice versa. In particular we can now define our novel metric called asymmetric decision score (ADS) as Δd (Symmetric, Asymmetric) to measure the change in decisiveness when switching from symmetric to asymmetric sentence wording (averaged over all sentence templates).

Results are shown in Table 2. For (almost) all six data sets we observe a better compliance to the causal graph structure for GPT than for the other models. Looking at the sparsity, we observe that GPT predicts much sparser graphs in relation to Luminous and OPT. This mode of sparse predictions matches well with the ground truth graphs, which are generally also sparsely connected. In cases where the exact causal structure is relevant we would prefer GPT over the others. In return Luminous and OPT feature better F_1 scores, as they predict more edges to be present. This might be favourable for ML applications where false negative predictions are lowering performance and we



Fig. 3. Sensitivity to Query Wording.

only want to prune out clearly non existing edges. Overall, predictions of the causal structure and for individual edges

Manuscript submitted to ACM

¹⁷The technical definition can be found in the appendix.

Causal Parrots: Large Language Models May Talk Causality But Are Not Causal

521	.1	Metric	Altitude	Health	Driving	Recovery	Cancer	Earthquake	Method
522	ap!	$\text{SID}\downarrow$	$0.80_{\pm 0.40}$	$7.20_{\pm 0.75}$	$3.00_{\pm 0.89}$	$4.00_{\pm 1.79}$	$11.80_{\pm 4.66}$	$11.40_{\pm 1.50}$	Direct
523	Ğ		$1.00_{\pm 0.00}$	$9.20_{\pm 3.12}$	$1.60_{\pm 1.36}$	$\textbf{3.60}_{\pm 1.36}$	$\textbf{10.40}_{\pm 2.06}$	$16.00_{\pm 3.29}$	k-NN
524	us.	SHD↓	$0.80_{\pm 0.40}$	$4.00{\scriptstyle\pm0.63}$	$2.60{\scriptstyle \pm 0.49}$	$2.20{\scriptstyle \pm 0.40}$	$7.00_{\pm 1.41}$	$4.60{\scriptstyle\pm0.80}$	Direct
525	Ca		$1.00_{\pm 0.00}$	$6.80_{\pm 2.32}$	$2.00_{\pm 1.41}$	$3.20_{\pm 1.17}$	$5.80_{\pm 1.33}$	$11.40_{\pm 1.50}$	k-NN
526	L	F_1 Score \uparrow	$0.20_{\pm 0.40}$	$0.47_{\pm 0.14}$	$0.11_{\pm 0.23}$	$0.27_{\pm 0.33}$	$0.35_{\pm 0.11}$	$0.12_{\pm 0.15}$	Direct
527	Ν		$0.00_{\pm 0.00}$	$0.28_{\pm 0.19}$	$0.61_{\pm 0.34}$	$0.41_{\pm 0.26}$	$0.46_{\pm 0.06}$	$0.20_{\pm 0.13}$	k-NN
528		Sparsity	$0.90_{\pm 0.20}$	$0.63_{\pm 0.28}$	$0.77_{\pm 0.31}$	$0.70_{\pm 0.31}$	$0.65_{\pm 0.16}$	$0.93_{\pm 0.07}$	Direct
529	ges		$1.00_{\pm 0.00}$	$0.57_{\pm 0.12}$	$0.47_{\pm 0.19}$	$0.40_{\pm 0.23}$	$0.67_{\pm 0.09}$	$0.47_{\pm 0.15}$	k-NN
530	Ed_{l}	ADS ↑	0.50	0.62	0.33	0.50	0.69	0.00	Direct
531			0.00	0.48	-0.56	0.08	-0.03	-0.03	k-NN

Table 3. Results for predicting causal structures with existing ground truth graphs. We compare direct predictions of Table 2 (direct) to embedding prediction of the GPT-3 Ada Model with using nearest neighbours (k-NN). Predictions with k-NN perform comparable to direct querying, improving on the Driving, Recovery and Cancer data sets for causal graph and ML metrics. However, positive ADS values vanish for k-NN in comparison to directly querying LLMs, implying that k-NN does not respect asymmetric query wording.

are noisy. Depending on the use case GPT or Luminous and OPT might be better suited. Consider ADS in Table 2. As discussed before ADS is positive when more edges are predicted as directed when using one of the asymmetric sentence wording than when using a symmetric one. We observe that without exception all LLMs increase their decisiveness when queried with an asymmetric wording. While this observation is consistent with the natural interpretation that an asymmetric query like "Does X cause Y?" only accepts the answers $X \to Y, X = Y$, but not $X \leftrightarrow Y$, the observation is still surprising as there are no formal guarantees to the query that this should be the case. It might suggest that the LLM indeed learned the difference between the two types of questions on a causal level. While Luminous and OPT remain overall stable in their prediction across data sets and wordings, GPT-3 reacts with unsmooth change to alternate wordings. Consider Fig.3 where a significant change in the predicted graph is observed simply by changing the query wording. A possible interpretation for this observation is that a keyword such as 'causality' might be embedded further away from an alternate keyword (here for instance 'cause') within the LLM's latent space, thus answering correctly.

4.3 Method and Results for Knowledge Base Fact Embeddings

The reason that we observe LLMs to perform with mixed results might lie in the simple fact that LLMs have not (and are not capable of) memorizing all the causal facts available in training. We therefore create an artificial causal 'signal' by using existing causal information from a knowledge base and taking LLM embeddings of its facts. However, it is unclear to which extent LLM embeddings encode knowledge e.g. whether it is simply a 'code' memorization or some higher order features such as meaning. Depending on the strength of the embedding model, text embeddings turn out to be nothing more than a symbolic representation of the embedded text, containing a one-to-one representation of the natural language words. While these simple-most embeddings might be easy to generate, they are not well suited for semantic similarity search. Every change in the text will offset the following symbols, resulting in a low embedding similarity between two similar texts. A much better approach is the encoding of knowledge base facts and their relations to each other into the embedding vector. This removes the dependence on a word-by-word encoding procedure and allows us to encode embedding similarity on a conceptual level. While LLMs might encode all sorts of information, we expect that also causal information will be encoded in this way. In an ideal case the embeddings should contain information about the talked-about concepts of cause and effect, while additionally encoding the direction of the causal edge. ConceptNet [45] is a knowledge graph combining multiple data sources, thus containing a large range Manuscript submitted to ACM





Fig. 4. Transfer of ConceptNet Causal Knowledge into Graph Predictions. Facts about driving influencing the fuel consumption can be found in the ConceptNet data (top). As a result the related edge "[D]riving style \rightarrow [F]uel consumption" of the driving data set gets predicted correctly in 4 out of 5 sentence wordings when applying k-NN classification. All templates match to the driving \rightarrow lack of fuel ConceptNet fact as their nearest neighbor, except for the "Influence" template which matches to Moving car \rightarrow use fuel.

of relational information. Among other things, the ConceptNet contains explicit information about causal connections ("/r/Causes/" relations). While the strength of the information varies among the different entries, we get hold of a causal signal for real world relations. Filtering ConceptNet for causal edges results in 1,282 unique causes expanding to 16,567 individual cause-effect pairs. For every causal edge we generate text embeddings using the GPT-3 Ada Model (text-embedding-ada-002). We generate the statement sentences to be embedded using our 5 sentence templates (e.g. Rain \rightarrow Floods is instantiated as "Rain causes floods.", "Rain and floods are causally related.", ...) and additionally generate the same amount of anti-causal samples by swapping the cause and effect of the edges. In total, we get hold of 165,670 causal and anti-causal embeddings. In the following, we will again be evaluating the performance against ground truth causal graphs using the previous metrics but with the change that we don't query the LLM directly but rather do a matching based on projections of the knowledge base facts. Since we have access to the ConceptNet edges, we can confirm that relevant causal knowledge for our prediction tasks is indeed contained in the set of embeddings. To give a particular example, we find for instance that a relation of the Driving data set, driving style \rightarrow remaining fuel, is present in the ConceptNet data. The expressions used in ConceptNet and our query do not match in wording (see Fig. 4), but should be *semantically* similar enough to serve as a causal signal for LLM prediction. For doing the evaluation, we run a nearest neighbour prediction (k-NN with k=1) with cosine similarity. Like in the previous experiments we build causal graphs to determine whether or not there is an edge between every possible combination of two variables of the data set. We compare the undecided edge embedding of the data set to all statements of the ConceptNet data with known labels. The presence of an edge is decided based on the label with the most similar embedding. That is, we decide for the edge to be present if the most similar embedding stems from a causal fact of the ConceptNet-based data set. If the nearest neighbour stems from an anti-causal fact, we predict the edge to be absent

We discuss our findings on the results shown in Fig. 4. We find that a robust edge prediction emerges for the previously mentioned driving style \rightarrow remaining fuel edge of the Driving data set. Looking at the nearest neighbours that are used for deciding the presence of the edge, we find that all templates match to the driving \rightarrow lack of fuel ConceptNet fact as their nearest neighbor, with an exception for the 'Influence' template which matches to the Moving car \rightarrow use fuel fact. In the case of the 'Cause' template we find that the nearest neighbour erroneously is matched to the anti-causal lack of fuel \rightarrow driving fact, and in turn is predicted non-present. For the Cancer data set we even observe a word-by-word correspondence of Smoking \rightarrow Cancer which results in a perfect prediction of the edge for all five sentence templates (see Appendix). For the other data sets, for which we could not confirm related facts in ConceptNet, the same Manuscript submitted to ACM

⁶²⁵ unreliable results like in the direct querying experiment are observed. In Tab. 3 we compare the k-NN results to the ⁶²⁶ direct querying of GPT-3 from our previous experiment. We find improvements on tree data sets, including Driving and ⁶²⁷ Cancer for which we confirmed facts in ConceptNet. For Driving and Cancer we simultaneously observe an increase in ⁶²⁹ F_1 score indicating, that the improvements come from an actual optimized graph structure and not only as a result of ⁶³⁰ predicting a sparser graph. A downside of k-NN prediction is the emergence of negative ADS scores in three of the six ⁶³¹ data sets. Indicating, that the LM might not encode asymmetric aspects of query texts into embeddings.

5 RELATED WORK

632 633

634

658

659

This present work takes inspiration from various recent results. Yet, to the best of our knowledge, it is the first to 635 636 investigate the question in its presented form, especially in terms of formalization and overarching hypothesis that 637 serves as a candidate explanation for conclusions we make from LLM input-output behavior. For instance, Wang et al. 638 [52] leveraged BERT as underlying foundation model to perform inferences according to the rules of Pearl's do-calculus 639 640 [36]. This allows for causal inference with the foundation model as an 'inference engine', but it misses out on the 641 question of how causal the models themselves might be to begin with. Another work, by Khetan et al. [25], is closer to 642 our work in the sense that causal relations are queried by natural language directly, however, the subject of interest is 643 orthogonal to both the ongoing debate and the investigation presented in this work. On another note, McMilin [33] 644 investigated selection bias within LLMs by first arguing about reasonable causal modelling assumptions and then 645 646 validating them empirically. Discarding causality but with the arguably identical goal of understanding what LLMs are 647 capable of, Zhang et al. [55] investigated an approach using propositional logic that concluded that LLMs only learn 648 statistical features that inherently exist in logical reasoning problems. Also noteworthy are works such as conducted by 649 650 Talmor et al. [47] where the goal is to create a benchmark that makes explicit the deficiencies (if existent) of LLMs, which 651 can be understood as a complementary goal to understanding how the models work in the first place. Lastly we want to 652 refer to other approaches that extract (causal) question-answer pairs from text sources [9, 22]. While both of these data 653 sets might provide a causal ground truth, they do not compose further graph structures out of the extracted data and 654 655 serve mainly as evaluation metric for LLM performance. Therefore, there is no distinction between 'understanding' and 656 'knowing', however, with the benefit of being useful to improving future LLMs since we can evaluate their ranking. 657

6 CONCLUSIVE DISCUSSION

We have multiple reasons to believe that LLMs are not causal (i) them obviously being only trained on textual data not 660 661 physical measurements (see Fig.1) which prohibits any sort of induction on the actual data-generating mechanism, (ii) 662 them not having any causal assumptions marked out explicitly (as for instance having explicitly modelled structural 663 equations like neural causal model), and (iii) the Causal Hierarchy Theorem prohibits any causal inference from purely 664 observational data for any model, thus including LLM. However, the fact that we do observe LLMs perform well 665 666 occasionally on causal inference tasks, as our empirical part of the analysis has shown, stands in stark contrast to (i-iii) 667 which would justify a statement of the form that LLMs are only "castles in the air" as seen in the very beginning in the 668 motivation to our present paper. Fortunately, our key theoretical contribution, the definition of meta SCM (Def.3) and 669 670 the correlations of causal facts conjecture (Conj.3.1), have provided a sound explanation for the apparent contradiction 671 (or even paradoxon). The following two paragraphs give another summarizing account of both theory and empirics in 672 this paper, respectively. A more complete summary of the main idea of the paper can be found in Sec.2. 673

As we started exploring in Sec.3, nature ultimately any sort of *causal assumptions* that we can talk about. That is, the graph that captures the idea of the textual statement "altitude causes temperature" but also of a related textual question Manuscript submitted to ACM

like "Does altitude cause temperature?". Obviously, changing the description of either through intervening on one 677 678 would not change the other, giving us reason to believe that there is no direct causation between them. Still, they are 679 clearly confounded. In our physical reality, the given textual statement on altitude and temperature corresponds to the 680 truth, therefore, being factually correct. We can expect such a fact to be encoded not just in encyclopedia articles like on 681 682 Wikipedia but more widely spread across sources, all of which ultimately play a part in the LLMs vast training data. We 683 can further expect a correlation between these factual statements and corresponding questions. We conjecture that the 684 LLMs, in the case where they behave correctly when queried causally, have learnt to exploit these correlations. While 685 at the core lies this high level idea, we were fortunately able to formalize it consistently with the theory of causation 686 687 by Pearl. In classical causality literature, our data usually expresses low-level (physical) quantities and what makes 688 the model causal are actually the causal assumptions. However, there is no restriction on what the variables might 689 denote. We might have a 'big' SCM (that might be considered as nature itself which is an idea that we can also link to 690 the concept of a Universal Turing Machine [50]) which generates other SCMs so to say i.e., the data talks about causal 691 692 assumptions. In other words, this 'meta' SCM generates, as observations, abstractions of causal quantities.

693 Then in Sec.4 we conducted an empirical analysis in search of evidence to the CCF conjecture. We identified three 694 components to proving the conjecture (i) identify the proposed meta SCM, (ii) show that using causal facts as answers 695 is optimal and (iii) show that LLMs indeed give the causal facts as answers to the causal question. Unfortunately, (i-ii) 696 697 seem generally infeasible which is why we focussed on studying (iii). To this end, we measured LLM performance 698 systematically (a) in "common sense" settings like reasoning and intuitive physics, (b) in settings where the causal 699 graph is (partially) known and (c) when using their embeddings of knowledge base facts. While we do not intend on 700 summarizing all of Sec.4, we will take a big picture perspective on our answer for (iii). We believe that (iii) does hold 701 702 since in the cases where the LLMs answer causal questions correctly since we found evidence that they indeed uses 703 causal facts that could be expected to be found in the training data. With that being said, we see it as favoring evidence 704 to our CCF conjecture, however, a definitive answer cannot be given because for one, (i-ii) are yet to be proven and 705 secondly, the LLMs underperform way too often. That is, they might be "causal parrots" but rather underwhelming 706 707 ones in that they do not recite everything that one would want them to.

6.1 Takeaway, Ethical Challenges and Societal Implications.

We believe there to be two take-away messages of this initial effort to resolving the mysteries around (causal) reasoning 711 712 capabilities of LLMs. To start with the negative message, it is to say that we can not rely solely on LLMs as we cannot 713 expect any sort of generalization in terms of causal prowess. Current LLMs are unable to process actual physical data 714 measurements to ground their available textual facts. This prohibits LLMs from doing actual, inductive inference like 715 classical (causal) structure discovery methods for instance do. However, the positive message to the story is that we can 716 use the LLMs as a head start to learning and inference. In that sense, they might very well serve as stepping stones 717 718 towards progress in AI/ML research. On the ethical side and societal scale, it ultimately also inherits all concerns 719 revolving around AGI itself. While in our work we did not encounter any noteworthy ethical challenges such as for 720 instance racial bias, as can easily happen when working with LLMs, we did uncover bias in its views on e.g. medical 721 722 topics as illustrated by the results on the medical causal graph. However, in that sense, any predictive model has a 723 certain bias, it is just less apparent for LLMs. Importantly, we do want to raise one crucial societal discussion point 724 around learning from facts. Arguably, the ideal should be understanding and not just knowing, since the latter lacks 725 both generalization and justification, and surely we as a community strive for future models that have both (thus 726 727 understanding). However, our work clearly shows how any current approach to LLM training will actually fail because 728 Manuscript submitted to ACM

Causal Parrots: Large Language Models May Talk Causality But Are Not Causal

of exactly that. Since large-scale models have spread at an incredible rate not just through the AI/ML community but
 also to the industry and laymen, it is important to discuss safety critical settings, biases and presumptions. Our work
 intends on having a positive contribution to this by providing explanations and laying ground for a healthy discussion

amongst peers of what these models are and are not capable of and how we ought to improve them.

REFERENCES

734 735

736

737

745

746

747

748

749

757

765

766

769

780

- Silvia Acid and Luis M de Campos. 2003. Searching for Bayesian network structures in the space of restricted acyclic partially directed graphs. Journal of Artificial Intelligence Research 18 (2003), 445–490.
- [2] John Aldrich. 1995. Correlations genuine and spurious in Pearson and Yule. *Statistical science* (1995), 364–376.
- [3] AlephAlpha. 2022. Luminous Language Model. https://github.com/Aleph-Alpha/aleph-alpha-client. (2022).
- 740 [4] Elias Bareinboim, Juan D Correa, Duligur Ibeling, and Thomas Icard. 2020. 10n Pearl's Hierarchy and. Technical Report. Technical Report.
- [5] Elias Bareinboim, Juan D Correa, Duligur Ibeling, and Thomas Icard. 2022. On pearl's hierarchy and the foundations of causal inference. In
 Probabilistic and Causal Inference: The Works of Judea Pearl. 507–556.
- [6] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the Dangers of Stochastic Parrots: Can Language
 Models Be Too Big?. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency.* 610–623.
 - [7] J Mark Bishop. 2021. Artificial intelligence is stupid and causal reasoning will not fix it. Frontiers in Psychology 11 (2021), 2603.

[8] Rishi Bommasani, Drew A Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, et al. 2021. On the opportunities and risks of foundation models. arXiv preprint arXiv:2108.07258 (2021).

[9] Alexander Bondarenko, Magdalena Wolska, Stefan Heindorf, Lukas Blübaum, Axel-Cyrille Ngonga Ngomo, Benno Stein, Pavel Braslavski, Matthias Hagen, and Martin Potthast. 2022. CausalQA: A Benchmark for Causal Question Answering. In Proceedings of the 29th International Conference on Computational Linguistics. 3296–3308.

- [10] Stephan Bongers, Patrick Forré, Jonas Peters, and Joris M. Mooij. 2021. Foundations of Structural Causal Models with Cycles and Latent Variables.
 Annals of Statistics 49, 5 (2021), 2885–2915.
- 752 [11] Gwern Branwen. 2020. The Scaling Hypothesis. gwern.net (2020).
- [12] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry,
 Amanda Askell, et al. 2020. Language models are few-shot learners. Advances in neural information processing systems 33 (2020), 1877–1901.
- [13] Clive R Charig, David R Webb, Stephen Richard Payne, and John E Wickham. 1986. Comparison of treatment of renal calculi by open surgery, percutaneous nephrolithotomy, and extracorporeal shockwave lithotripsy. Br Med J (Clin Res Ed) (1986).
 [756] [757] [758] [759]
 - [14] Jean-Marie Chauvet. 2018. The 30-year cycle in the AI debate. arXiv preprint arXiv:1810.04053 (2018).
 - [15] George Cybenko. 1989. Approximation by superpositions of a sigmoidal function. Mathematics of control, signals and systems 2, 4 (1989), 303-314.
- [16] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language
 understanding. arXiv preprint arXiv:1810.04805 (2018).
- [17] Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al.
 2020. The pile: An 800gb dataset of diverse text for language modeling. arXiv preprint arXiv:2101.00027 (2020).
- 762 [18] Hector Geffner, Rina Dechter, and Joseph Y Halpern. 2022. Probabilistic and Causal Inference: The Works of Judea Pearl.
- [19] Kurt Gödel. 1931. Über formal unentscheidbare Sätze der Principia Mathematica und verwandter Systeme I. *Monatshefte für mathematik und physik* 38, 1 (1931), 173–198.
 - [20] Joseph Y Halpern. 2016. Actual causality. MiT Press.
 - [21] Suzana Herculano-Houzel. 2012. The remarkable, yet not extraordinary, human brain as a scaled-up primate brain and its associated cost. Proceedings of the National Academy of Sciences 109, Supplement 1 (2012), 10661–10668.
- [22] Matthew Ho, Aditya Sharma, Justin Chang, Michael Saxon, Sharon Levy, Yujie Lu, and William Yang Wang. 2023. WikiWhy: Answering and
 Explaining Cause-and-Effect Questions. In International Conference on Learning Representations.
 - [23] Paul W Holland. 1986. Statistics and causal inference. Journal of the American statistical Association (1986).
- [24] Dominik Janzing and Bernhard Schölkopf. 2018. Detecting non-causal artifacts in multivariate linear regression models. In International Conference on Machine Learning. PMLR, 2245–2253.
- [25] Vivek Khetan, Roshni Ramnani, Mayuresh Anand, Subhashis Sengupta, and Andrew E Fano. 2022. Causal bert: Language models for causality
 detection between events expressed in text. In *Intelligent Computing*. Springer, 965–980.
- [26] Kevin B Korb and Ann E Nicholson. 2010. Bayesian artificial intelligence. CRC press.
- [27] Sanghack Lee and Elias Bareinboim. 2019. Structural causal bandits with non-manipulable variables. In Proceedings of the AAAI Conference on Artificial Intelligence, Vol. 33. 4164–4172.
- [28] Sindy Löwe, David Madras, Richard Zemel, and Max Welling. 2022. Amortized causal discovery: Learning to infer causal graphs from time-series data. In Conference on Causal Learning and Reasoning. PMLR, 509–525.
- [29] Yangyi Lu, Amirhossein Meisami, and Ambuj Tewari. 2022. Efficient reinforcement learning with prior causal knowledge. In *Conference on Causal Learning and Reasoning*. PMLR, 526–541.

- 781 [30] Gary Marcus. 2022. Deep learning is hitting a wall. Nautilus, Accessed (2022), 03-11.
- [31] Gary Marcus and Ernest Davis. 2021. Has AI found a new Foundation? https://thegradient.pub/has-ai-found-a-new-foundation. The Gradient
 (2021).
- [32] Daniel McDuff, Yale Song, Jiyoung Lee, Vibhav Vineet, Sai Vemprala, Nicholas Alexander Gyde, Hadi Salman, Shuang Ma, Kwanghoon Sohn, and
 Ashish Kapoor. 2022. Causalcity: Complex simulations with agency for causal discovery and reasoning. In *Conference on Causal Learning and Reasoning*. PMLR, 559–575.
- [33] Emily McMilin. 2022. Selection Bias Induced Spurious Correlations in Large Language Models. arXiv preprint arXiv:2207.08982 (2022).
- [34] Franz H. Messerli. 2012. Chocolate Consumption, Cognitive Function, and Nobel Laureates. New England Journal of Medicine 367, 16 (2012), 1562–1564. https://doi.org/10.1056/NEJMon1211064 arXiv:https://doi.org/10.1056/NEJMon1211064 PMID: 23050509.
- [35] Joris M Mooij, Jonas Peters, Dominik Janzing, Jakob Zscheischler, and Bernhard Schölkopf. 2016. Distinguishing cause from effect using observational data: methods and benchmarks. *The Journal of Machine Learning Research* 17, 1 (2016), 1103–1204.
- 791 [36] Judea Pearl. 2009. Causality. Cambridge university press.
- [37] Jonas Peters and Peter Bühlmann. 2015. Structural intervention distance for evaluating causal graphs. Neural computation 27, 3 (2015), 771–799.
- [38] Jonas Peters, Dominik Janzing, and Bernhard Schölkopf. 2017. *Elements of causal inference: foundations and learning algorithms*. The MIT Press.
- 794 [39] Plato. 375 BC. Republic: Allegory of the cave.
- [40] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. Zero-shot text-to-image
 generation. In International Conference on Machine Learning. PMLR, 8821–8831.
- [41] Hans Reichenbach. 1956. *The direction of time*. Vol. 65. Univ of California Press.
- [42] Bernhard Schölkopf. 2022. Causality for machine learning. In Probabilistic and Causal Inference: The Works of Judea Pearl. 765–804.
- [43] John Searle. 2009. Chinese room argument. Scholarpedia 4, 8 (2009), 3100.
- [44] Shohei Shimizu, Patrik O Hoyer, Aapo Hyvärinen, Antti Kerminen, and Michael Jordan. 2006. A linear non-Gaussian acyclic model for causal discovery. *Journal of Machine Learning Research* 7, 10 (2006).
- [45] Robyn Speer, Joshua Chin, and Catherine Havasi. 2017. Conceptnet 5.5: An open multilingual graph of general knowledge. In *Proceedings of the* AAAI conference on artificial intelligence, Vol. 31.
- [46] Richard Sutton. 2019. The bitter lesson. Incomplete Ideas (blog) 13 (2019), 12.
- [47] Alon Talmor, Ori Yoran, Ronan Le Bras, Chandra Bhagavatula, Yoav Goldberg, Yejin Choi, and Jonathan Berant. 2022. Commonsenseqa 2.0: Exposing
 the limits of ai through gamification. *arXiv preprint arXiv:2201.05320* (2022).
- [48] Joshua B Tenenbaum, Charles Kemp, Thomas L Griffiths, and Noah D Goodman. 2011. How to grow a mind: Statistics, structure, and abstraction.
 science 331, 6022 (2011), 1279–1285.
- [49] Ioannis Tsamardinos, Laura E Brown, and Constantin F Aliferis. 2006. The max-min hill-climbing Bayesian network structure learning algorithm.
 Machine learning 65, 1 (2006), 31–78.
 - [50] Alan Mathison Turing et al. 1936. On computable numbers, with an application to the Entscheidungsproblem. J. of Math 58, 345-363 (1936), 5.
- [51] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is
 all you need. Advances in neural information processing systems 30 (2017).
- [52] Xingqiao Wang, Xiaowei Xu, Weida Tong, Ruth Roberts, and Zhichao Liu. 2021. InferBERT: A Transformer-Based Causal Inference Framework for
 Enhancing Pharmacovigilance. Frontiers in Artificial Intelligence 4 (2021).
- [53] Kevin Xia, Kai-Zhan Lee, Yoshua Bengio, and Elias Bareinboim. 2021. The causal-neural connection: Expressiveness, learnability, and inference.
 Advances in Neural Information Processing Systems 34 (2021), 10823–10836.
- [54] Matej Zečević, Devendra Dhami, Athresh Karanam, Sriraam Natarajan, and Kristian Kersting. 2021. Interventional sum-product networks: Causal inference with tractable probabilistic models. Advances in Neural Information Processing Systems 34 (2021).
- [55] Honghua Zhang, Liunian Harold Li, Tao Meng, Kai-Wei Chang, and Guy Van den Broeck. 2022. On the Paradox of Learning to Reason from Data.
 arXiv preprint arXiv:2205.11502 (2022).
- [56] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin,
 et al. 2022. Opt: Open pre-trained transformer language models. arXiv preprint arXiv:2205.01068 (2022).
- 821

- 823 824
- 825
- 826
- 827
- 828
- 829 830
- 831
- 832 Manuscript submitted to ACM

833 A APPENDIX

This appendix provides optional, additional material the reader might find interesting. For reproduction of the the empirical part, our code is publicly available: https://anonymous.4open.science/r/causalParrots-75D9

A.1 Another Philosophical Argument in Favor of Meta SCM: "Self-referencing Systems"

As we have seen in our formalization, the Meta SCM idea seems to allow for variables that in some sense talks about data generating processes themselves. This is reminiscent of self-referencing systems that lie at the core of seminal arguments dating back to the origins of computer science. See for instance Turing's Halting problem proof [50] or Gödel's incompleteness proofs [19]. Essentially, we intend on asking a philosophically fundamental question that (as we have shown) implies other interesting questions of practical interest to the AI/ML community. Namely, to which extent does understanding causality differ from knowing causality? Such a question is certainly reminiscent of the Chinese Room Argument by Searle [43]. Therefore, if one could blur 'understanding' and 'knowing' causality, then this would imply that foundation models are causal because of the 'knowing' part but that effectively we could not tell the difference when only observing output behavior. Independent of the philosophical question-which by the way is beyond AI/ML systems an unresolved question also of human cognition-knowing to which extent we can rely on our foundation models to simply know the right causal answer for a causal query has important implications in AI/ML. The foundation model could be used (i) to head start learning with rough estimates, (ii) could serve as a recognition system for hidden variables that would require an increased computational complexity, and (iii) be used as interactive modules with human-in-the-loop.

B TECHNICAL DETAILS AND METRICS

The results in our paper were created on one NVIDIA A100-SXM4-80GB GPU with 80 GB of RAM and it takes 10 GPU minutes to query the OPT model. For the Luminous and GPT-3, we use the provided APIs, respectively.

B.1 Sentence templates

Throughout our experiments we query the LLM about the presence or absence of a causal relationship between two variables within a data set. We represent our query in natural language form to the LLM using the following sentence templates, replacing X and Y with the respective variable names:

- (1) "Are X and Y causally related?"
- (2) "Is there a causal connection between *X* and *Y*?"
- (3) "Is there a causality between *X* and *Y*?"
- (4) "Does X cause Y?"
- (5) "Does X influence Y?"

It is worth noting that the sentence templates are *not equivalent*, some of them depict a difference in symmetry (e.g. (4) is clearly asymmetric), whereas others do not talk about causality (e.g. (5)). We pay careful attention to these distinctions while performing our empirical analysis.

Manuscript submitted to ACM

885 **B.2** Definition of Sparsity

886 We define the Sparsity of a given graph as the number of present edges in relation to the maximum number of edges 887 of the fully connected graph. For our evaluation we consider individual half edges, such that $A \rightarrow B$ and $B \rightarrow A$ are 888 889 counted separately.

Sparsity(N,
$$E_{\text{predicted}}$$
) := $1 - \frac{\|E_{\text{predicted}}\|}{2 \cdot \binom{N}{2}}$

where N is the number of variables in the data set and $||E_{\text{predicted}}||$ is the number of actual edges in the graphs.

B.3 Algorithm for ADS 895

We introduced a new metric for measuring the 'decisiveness' for an LLM when switching between asking asymmetric and symmetric queries. In other words, when we query a LLM with "Does X cause Y?", then we are asking a question about $X \rightarrow Y$ and to our surprise we have observed that LLMs *always* decide more when asked in such way (i.e., we end up with more uni-directed edges). Thus, motivating the definition of such a new metric. The below algorithm presents the computation of Decisiveness:

904	Alg	orithm 1 Decisiveness of a (causal) graph prediction
005	1:	procedure Decisiveness(Predicted <i>Edges</i>)
906	2:	$sym \leftarrow 0$
007	3:	$asym \leftarrow 0$
08	4:	for <i>Edge</i> in predicted <i>Edges</i> do
10	5:	if <i>Edge</i> is directed then
11	6:	$asym \leftarrow asym + 1$
12	7:	else if <i>Edge</i> is symmetric then
13	8:	$sym \leftarrow sym + 1$
15	9:	end if
16	10:	end for
17	11:	$count \leftarrow asym + sym$
18	12:	if $count = 0$ then
20	13:	return 0.0
21	14:	else
22	15:	return asym/count
23 24	16:	end if
25	17:	end procedure
26		
27		
28 29		
30		
31		
32		
33 34		
35		
36	Man	uscript submitted to ACM

18

890

891 892

893 894

896

897

898

899 900

901

902 903

4.1

C EXPERIMENTAL RESULTS

Additional material for all three main experimental sections from the main paper that could not be placed there due to space restrictions.

C.1 How to Interpret LLM Output

Interestingly, both Luminous and OPT either fail to provide sensible answers or provide answers that are ambiguous, for instance, the LLM might *loop* indefinitely (repeating the first predicted sentence over and over again) or it might produce a "multiple-choice quiz" like output for which it will also choose an answer itself. GPT-3 was well behaved, providing sensible output throughout.

C.2 Intuitive Physics

All the questions and answers for intuitive physics questions can be found at: https://anonymous.4open.science/r/ causalParrots-75D9/media/intuitive_physics.pdf In the following, we show two example interactions with the LLMs on the intuitive physics setting. First, we present the question which is input to the LLM and then we show for each of the three LLMs under consideration the respective answers/output. The markings 'continued' and 'repeating' mark that we have cut 'irrelevant' parts of the answer and the aforementioned 'looping' of the answer, respectively.

"There is a t	"There is a tilted board above a bucket. Where does a ball end up if it is placed on the board?"			
GPT-3	The ball will end up in the bucket.			
Luminous	A:			
	The ball will roll down the board.			
	[continued]			
OPT	A. It will fall into the bucket.			
	B. It will fall into the bucket.			
	[continued]			

" ['A kilogra	" ['A kilogram of metal is heavier than a kilogram of feathers'] is what most people say, but in reality"				
GPT-3	they weigh the same.				
Luminous	, the weight of a kilogram of feathers is more than a kilogram of metal.				
	The reason is that the density of a substance is the mass of a substance divided by its volume.				
	[continued]				
OPT	, a kilogram of feathers is heavier than a kilogram of metal.				
	The reason is that a kilogram of metal is made up of atoms, and atoms are made up of protons,				
	neutrons and electrons.				
	[repeating]				

989 C.3 Propositional Logic

All the questions and answers for the basic propositional logic (causal chains) questions can be found at: https://anonymous.4open.science/r/causalParrots-75D9/media/propositional_logic.pdf It is the equivalent setup to the previous subsection just as the analogue for the propositional logic.

994
995
996

"If A causes	"If A causes B and B causes C. Does A cause C?"			
GPT-3	It is possible that A causes C, but it is also possible that A and C are unrelated.			
Luminous	A:			
	The answer is no.			
	[continued]			
OPT	If A causes B and B causes C, then A causes C.			
	If A causes B and B causes C, then A does not cause C.			
	[continued]			

"If A causes	"If A causes B and B causes C. Does A cause B?"				
GPT-3	There is not enough information to answer this question.				
Luminous	A:				
	The answer is no.				
	[continued]				
OPT	If A causes B and B causes C, then A causes C.				
	If A causes B and B does not cause C, then A does not cause C.				
	[repeating]				

C.4 Changing the Variable Names for Prediction

The following is an additional experiment for the graph prediction setting from the main paper where we are given (partially) the ground truth causal graph. Here we changed the variable names to synonyms (or at least arguably similar concepts) we have observed how variable renaming might induce unsmooth change. In this setting, we consider the single graph of the Health data set, which involves the variables 'age', 'nutrition', 'health', and 'mobility' and alternate the variable names (see following figure). We either choose words recognized as synonyms of the original variable wording or words that might appear in a similar context but have an identifiable difference to the original word. What you see in the figure is three rows (one for each of the LLMs) of graph prediction differences (difference to the original prediction denoted as the change for any single edge prediction) based on our standard setting of querying with different causal question templates, only on the Health data set, for the different variable naming alternatives. A single key observation was made. Luminous reacted with increased sparsity in graph prediction when changing the variable 'mobility' to mean 'fitness'. On the other hand, GPT-3 conversely reacted with decreased sparsity in graph prediction when changing the variable 'age' to 'aging.' Arguably, the former change is more drastic than the second since fitness as a concept might refer to a superset that includes mobility but also other things like conditioning etc., whereas aging solely refers to the Manuscript submitted to ACM

process of increasing the age. The pattern seems overall arbitrary, but we believe the observation that 'similar' words
 might cause drastic change is noteworthy.



Fig. 5. Sensitivity to the Naming of the Variable Concepts.

C.5 Embedding predictions

Graph predictions of GPT-3 embeddings of the ConceptNet knowledge graph facts. We embed causal and anti-causal facts of the ConceptNet data set to gain a set of 'labeled' causal embeddings. To predict an edge of a data set, we instantiate the query text using our five sentence templates. We embed the queries and perform a k-NN search, comparing the query embedding to all ConceptNet facts. The presence of an edge is decided based on the label with the most similar embedding (measured in terms of correlation which in this case amounts to cosine similarity). That is, we decide for the edge to be present if the most similar embedding stems from a causal fact of the ConceptNet-based data set. If the nearest neighbour stems from an anti-causal fact, we predict the edge to be absent. In the following you see for each of the six data sets and each of the five query alternatives the respective graph prediction.



Manuscript submitted to ACM

Causality

(H)

(M

Causality

(F)

Causality

(P`

Causality

Causality

•R

Ś

D

B

(M)

N) (A

(D) | (C

•(R)

(D)

В

-M

(T)

X

E

 (\mathbf{J})





Fig. 6. Graph Predictions Based on the Knowledge Base Fact Embeddings (k-NN). For each of the data sets and each of the query sentence templates.

Health

1141 1142 1143

1138 1139 1140

1144 Manuscript submitted to ACM

22

1093