Towards *Computational Comprehension*: A Non-Anthropocentric Framework for Evaluating LLM Understanding

Anonymous ACL submission

Abstract

This position paper introduces and motivates Computational Comprehension as an alternative, non-anthropocentric approach to assessing how Large Language Models (LLMs) handle knowledge. Unlike standard benchmarks, which often reward models for surface-level accuracy but shed little light on deeper conceptual understanding, Computational Comprehension directs attention to the model's internal processes. Specifically, we focus on whether certain neurons or sub-networks remain consistently activated across reformulations and contextual shifts of the same underlying concept. We outline a framework that tests a model's ability to preserve conceptual invariance under various input transformations, then observes how targeted ablations of relevant sub-networks affect performance. By gauging these internal, concept-related responses rather than relying solely on external metrics, we obtain more finegrained insights into a model's capacity to internalize, manipulate, and robustly apply conceptual knowledge. We also propose integrating such analysis into systematic experiments, showing how subtle tweaks to task prompts or data can reveal whether a model is genuinely concept-driven or merely parroting surface correlations. Through Computational Comprehension, we encourage researchers, engineers, and theorists to adopt a deeper, more transparent mode of evaluation-one that foregrounds internal conceptual grounding over score-centric arms races in pursuit of ever-higher benchmark numbers

1 Introduction

011

013

018

Designing and evaluating Large Language Models (LLMs) has long relied on human-centric notions of *understanding*, shaped by both philosophical definitions and our intuitive experience as humans. Traditionally, the term *understanding* invokes a cognitive process in which a person internalizes and applies concepts to model an event, context, or a piece of information (Bereiter, 2002). However, this anthropocentric framework creates immediate challenges for modern AI systems: ascribing human-style mental states or semantic understanding to models poses possibility of conflating symbol manipulation with genuine comprehension (Bennett, 2023). Recent arguments-including the Chinese Room thought experiment (Searle, 1980), the Mary's Room scenario (Jackson, 1986), and the Stochastic Parrot critique (Bender et al., 2021)-have raised doubts as to whether text-based LLMs can possess any more than surface-level patterns. Indeed, such arguments underscore a deeper question: in what sense can an artificial system be said to understand at all?

044

045

046

047

051

055

058

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

078

081

We propose a new term, **Computational Comprehension**, to distinguish the kind of understanding that may be attained by LLMs and other advanced neural architectures. Unlike human cognition, which is heavily shaped by embodiment, sensory modalities, and evolution (Foglia and Wilson, 2013), Computational Comprehension emphasizes the unique capacity of a model to form, store, and manipulate its own internal representations (or features) in a manner that enables consistent generalization and coherent behavior.

In doing so, we acknowledge both the findings of *mechanistic interpretability* research—which highlights emergent circuits that correlate with conceptlike features (Olah et al., 2020)—and the doubt that arises when we observe these models hallucinating, making errors, or merely pattern-matching pre-trained data. By framing an alternative, *nonanthropocentric* construct, our notion of Computational Comprehension aims to capture both the potential and the limitations of such systems without presupposing human-like consciousness or qualia.

A critical motivation for introducing and clarifying this term is the evolving crisis in LLM benchmarks. Many of today's LLMs not only surpass the average human in certain standardized tests but are rapidly approaching near-perfect performance on tasks designed by humans (Strachan et al., 2024). As a consequence, human score on benchmarks no longer provides a sufficient yardstick (Hern ´ andez-Orallo, 2020). If two models both exceed human excellence (near-100% accuracy), how should one meaningfully evaluate their capabilities? We argue that simply escalating benchmark difficulty is not a robust long-term solution. Instead, we must adopt an orthogonal, conceptual axis of evaluation: a model's degree of Computational Comprehension.

086

087

090

094

100

102

103

104

105

106

107

108

110

111

112

113

114

115

116

117

118

119

120

121

122

123

By investigating whether models truly capture, adapt, and integrate their internal features—as opposed to *merely* reproducing text patterns reviewed in training, we gain a more nuanced, explanatory picture of how they handle knowledge. Examining Computational Comprehension thus becomes a crucial research direction: it reframes the debate on LLMs' understanding, guides interpretability studies, and opens new vistas for building an evaluation framework that reflect real conceptual robustness.

In this paper, we elaborate upon the philosophical underpinnings that motivate this new terminology and propose initial pathways to measure Computational Comprehension in practice. By doing so, we seek to foster a framework that meaningfully captures the ways in which next-generation models process concepts, even if that process differs radically from human cognition. Our goal is not merely to defend a novel philosophical position but to engage practitioners and theorists alike in rethinking how we might define, test, and refine model-based comprehension without anchoring it to a solely anthropocentric conceptual lineage. Throughout this paper, we speak of model understanding and Computational Comprehension in the specific context of text-based LLMs.

2 Philosophical & Cognitive Foundations

In this section, we lay out the philosophical and 124 cognitive cornerstones that shape our understand-125 ing of understanding itself. By examining long-126 standing debates about human cognition, we high-127 light the anthropocentric assumptions often embed-128 ded in discussions of AI (Millière and Rathkopf, 129 130 2024). We begin by introducing how the human basis of understanding informs (and potentially bi-131 ases) our interpretation of large language model ca-132 pabilities, then turn to three well-known arguments 133 that each challenge whether symbolic proficiency 134

can ever equate to genuine comprehension.

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

159

160

161

162

163

164

165

166

167

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

2.1 The Human Basis of Understanding

The concept of *understanding* has been explored by philosophers and cognitive scientists for centuries, commonly framed as a process wherein an individual internalizes and applies concepts to interpret or predict phenomena in the world (Bereiter, 2002). However, these perspectives often adopt a humancentric perspective, or anthropocentrism, which positions human beings at the center of consideration. In this viewpoint, everything from philosophical inquiries to recent developments, such as LLMs, are often evaluated through human standards, reflecting the belief that human knowledge, reasoning, and consciousness are the ultimate criteria by which we measure understanding.

Yet the question remains whether this anthropocentric concepts translates readily to artificial architectures. These architectures may excel at pattern recognition without replicating *human subjective dimensions* (consciousness, experiences, etc.), which philosophers often regard essential to genuine comprehension (Foglia and Wilson, 2013). This raises the question *Can LLMs' internal parameters actually constitute features analogous to human concepts, or are they simply finely tuned statistical representations?*

2.2 Central Challenges: Three Arguments

To answer this question, we begin by examining *Searle's Chinese Room* (Searle, 1980), *Mary's Room* (Jackson, 1986), and the *Stochastic Parrot* critique (Bender et al., 2021), each aiming to highlight a distinct aspect of the gap between symbolic prowess and what human-centric notions would consider true comprehension.

Searle's Chinese Room

John Searle's Chinese Room thought experiment argues that the manipulation of symbols according to formal rules—yielding outputs indistinguishable from those of a fluent speaker—does not constitute genuine understanding. In this scenario, a person who speaks only English follows instructions to map Chinese input symbols to output symbols, appearing to understand Chinese, yet lacks any real semantic comprehension. This mirrors concerns in AI that models may generate appropriate responses without internal comprehension, effectively acting as sophisticated rulebooks where substantial parameter weights guide permissible patterns with184out innate semantic grounding. Some argue that185modern neural networks might develop emergent186properties transcending simple syntactic manipu-187lation (Wei et al., 2022), but whether these suffice188for true semantic understanding remains highly de-189bated, making Searle's thought experiment crucial190to current discussions about AI cognition.

Mary's Room

191

215

In Frank Jackson's narrative, Mary is a color sci-192 entist with complete theoretical knowledge about 193 color perception, yet she has only experienced a 194 grayscale world. Upon actually seeing red, she 195 gains new insights-qualitative knowledge about 196 color that no enumeration of facts could convey. 197 Initially challenging materialism in the philosophy of mind, this scenario also highlights the sig-199 nificance of subjective experience, beyond mere propositional facts, in understanding. For AI, the Mary's Room thought experiment underscores how purely symbolic or descriptive mastery might fail to capture certain dimensions of understanding, particularly those reliant on firsthand, qualitative experience. Since most LLMs are trained on linguis-206 tic data alone (Collier et al., 2022), critics argue 207 these models are inherently limited to factual or 208 correlational knowledge without the capacity for 209 perceptual immersion. For those who see a tight bond between experience and comprehension, the 211 thought experiment provides a powerful rejoinder 212 to claims that advanced text models truly know or 213 understand phenomena like color. 214

Stochastic Parrot Critique

Bender and colleagues (Bender et al., 2021) ar-216 gue that LLMs generate text similar to meaningful 217 human discourse but do so primarily through prob-218 219 abilistic rearrangement of learned patterns. These systems, as a result, assemble outputs in a manner that exhibits impressive fluency without necessarily having any underlying conceptual or intentional structure. From this perspective, an AI's knowledge is a patchwork drawn from extensive training data, 224 appearing coherent because it overlays existing linguistic tropes rather than genuinely grasping the subject matter. This critique warns against attributing true *understanding* simply because a machine's output aligns with human expectations, especially as LLMs excel at tasks traditionally associated with human cognition. Proponents argue that sufficiently complex systems might develop emergent properties beyond mere pattern manipulation (Wei et al., 233

2022). Nonetheless, we face the dilemma of discerning whether we are witnessing a qualitatively234novel form of cognition or merely the expansion of236surface-level pattern recognition.237

238

240

241

242

243

244

245

246

247

248

249

250

251

252

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

268

269

270

271

272

273

274

275

276

277

278

279

3 Hints of Artificial Comprehension

Mechanistic Interpretability

Mechanistic interpretability aims to open the *black* box of large-scale neural networks by uncovering how internal components-layers, neurons, attention heads, or circuits-instantiate concept-specific computations (Bereska and Gavves, 2024). Various interpretability techniques reveal a model's internal "concept-like" structures, or features (Templeton et al., 2024; Gao et al., 2024). Circuit tracing with probing classifiers identifies neurons that consistently fire for certain inputs (Olah et al., 2020), while feature visualization and activation patching highlight the tokens or activations most responsible for concept processing (Heimersheim and Nanda, 2024; Meng et al., 2022). Ablation studies further show that disabling suspect neurons degrades performance on targeted tasks (Yu and Ananiadou, 2024b; Zhang et al., 2024; Yu and Ananiadou, 2024a). By these means, we can often pinpoint how a model implicitly organizes its representation of knowledge. Empirical work demonstrates that removing these concept-circuits collapses functionality (Olah et al., 2020), suggesting dedicated internal pathways for distinct functions.

Emergent Phenomena

Grokking (Power et al., 2022) and prompt engineering such as Pause Tokens (Goyal et al., 2023) both exemplify intriguing emergent phenomena that hint at latent conceptual capacities. In grokking, models may plateau for many epochs before abruptly converging to near-perfect performance, suggesting a phase transition from rote memorization to a more generalizable, concept-based representation. Likewise, inserting seemingly trivial instructions or pause tokens can dramatically alter a model's output, apparently switching on previously dormant internal circuits. Mechanistic analyses of these effects indicate that large language models can harbor latent conceptual structures (Nanda et al., 2023), which remain quiescent until triggered by just the right input cues.

327

328

329

Next Step: Computational Comprehension

281

284

290

291

296

301

306

310

311

312

313

314

315

317

319

321

322

323

326

Mechanistic interpretability and emergent phenomena each shed light on a model's capacity for true Computational Comprehension. From circuit tracing and ablation studies, we learn how knowledge may be internally structured: identifying subnetworks that directly handle abstract concepts suggests more than mere pattern memorization. Likewise, phenomena like grokking and the dramatic effects of pause tokens highlight a network's latent dispositions, revealing how a seemingly small change in prompts or training conditions can trigger a deeper, more coherent conceptual state. Together, these findings go beyond measuring whether a model performs well on benchmarks: they offer insight into how robustly and systematically the model organizes its knowledge, thereby hinting at (or refuting) the presence of enduring, conceptlike structures that distinguish superficial patterning from genuine comprehension.

4 Towards Computational Comprehension

Recent LLMs exhibit near human performance, while arguments from section 2 emphasize the inability of LLMs to acquire the experiential or conscious dimensions typically associated with genuine comprehension. We introduce Computational Comprehension as a middle ground approach: an explicitly *non*-anthropocentric notion of what it might mean for a machine to *understand* while remaining faithful to the underlying computational mechanisms and constraints.

4.1 Defining Computational Comprehension

We propose that Computational Comprehension refers to the capacity of a computational model to:

- Internalize and represent a stable concept structure through its internal parameters or sub-networks, rather than mere memorization of specific input–output pairs.
- Exhibit systematic responses across diverse contexts that test multiple dimensions of a given concept, demonstrating an ability to generalize.
- Manipulate these internal concept-like representations in a way that allows for novel combinations, inferences, or transformation that align with the concept's formal or functional properties.

4.2 Key Dimensions of Computational Comprehension

Alongside the core definition, we propose several dimensions that can help refine how one might observe or measure Computational Comprehension in practice:

Generality & Robustness Possessing generality and being robust allows a model to effectively apply its understanding of concepts across various tasks and contexts, demonstrating that its knowledge goes beyond memorization. For example, a system trained for addition should correctly solve novel sums presented in different formats, including newly structured inputs or unfamiliar storybased settings, and it should maintain accuracy when problems are paraphrased or slightly altered, rather than failing due to superficial modifications. Together, these attributes emphasize the depth of conceptual understanding necessary for effective machine comprehension.

Sub-Network Localization and Interpretability Identifying neurons, attention heads, or circuits that consistently activate for a specific concept offers deeper insight into whether the concept is truly internalized. Coherent sub-networks that prove indispensable across tasks relating to the concept point to a stable, localized representation.

Compositional Integration If the model can manipulate concept A and concept B, does it also reliably handle combinations of A and B (e.g., nested constructs, intersection of properties), a hallmark of truly concept-based reasoning?

This characterization steers away from ascribing any notion of *phenomenal consciousness* or *qualia* to the model. Instead, *comprehension* here is pinned to the *use* of internally stored features that enable conceptual coherence. By being explicit about what these features (and their use) look like in computational terms, we aim to avoid unnecessary anthropocentrism and keep the focus on *how* the LLMs manipulate internal concepts into consistent outputs.

5 Benchmarking, Limitations, and the Need for a New Paradigm

As large language models (LLMs) have rapidly improved, benchmark performance has been the primary metric for gauging progress. While advancements in tasks like translation and question answering are impressive, it remains unclear whether high scores on these metrics indicate any deeper form of *understanding*(Mondorf and Plank, 2024), or Computational Comprehension. We examine the limitations of benchmark-driven evaluations and argue for a new paradigm, focused on concept-based perspectives.

376

377

385

397

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

5.1 Historical Context of Benchmarks and Accelerating Model Capabilities

Traditional benchmarks such as GLUE (Wang, 2018), SuperGLUE (Wang et al., 2019), SQuAD (Rajpurkar, 2016), MMLU (Hendrycks et al., 2020), and Big-Bench (Srivastava et al., 2022) have become standard metrics for evaluating progress in natural language processing (Fourrier et al., 2024), while modern LLMs achieve scores that rival or surpass average human baselines. While these advancements demonstrate the models' ability to capture complex statistical regularities, a concerning pattern emerges: once a benchmark gains recognition, models are optimized to reach near the theoretical ceiling (Zhou et al., 2023). This leads us to question what these scores truly signify: do they reflect a meaningful *understanding* of the underlying principles, or are they merely the result of memorizing patterns and exploiting correlations?

5.2 Why High Benchmark Scores May Not Imply Real Comprehension

A major critique of high benchmark scores is that models can achieve strong performances through superficial pattern matching and memorization rather than true understanding (Bender et al., 2021). For instance, when a dataset establishes consistent associations between specific words and outcomes, a model may learn to recognize these patterns akin to a *lookup table*, rather than forming robust, concept-level generalizations. This phenomenon can lead a model to excel on a particular dataset yet struggle with slight rephrasings or different contexts of the same task. Consequently, this raises the question of whether such high scores truly reflect Computational Comprehension.

5.3 Human-Centric Benchmarking Dilemma Beyond Human-Level Evaluation

Many of the most widely recognized benchmarks
were either explicitly designed to capture difficulties humans face or implicitly assumed that human
performance provides an upper bound (Zhou et al.,
2023). Now that models are meeting or exceeding

human baselines, a paradox arises: if both Model A and Model B surpass average human performance, how do we evaluate which model is better in a conceptual sense? Traditional metrics like accuracy or F1 score, lose discriminative power once performance clusters near 100%. 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

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

Eschewing an Endless Benchmark Arms Race

An intuitive reaction might be to escalate benchmark difficulty (Arora et al., 2023), creating new sets of even more intricate questions that push models further. However, an endless cycle of difficult benchmark creation risks feeding back into a purely performance-driven arms race. Such an approach may continue yielding impressive numeric gains but could neglect how models handle concepts and whether they truly internalize them or simply scale up memorization. Unless we adopt new perspectives, the gap between high test accuracy and genuine conceptual capacity may widen.

5.4 Criteria for a New Paradigm

Orthogonal to Difficulty

A key realization is that we do not need to simply *increase* the difficulty of tasks. Instead, we must develop *orthogonal* criteria that evaluate how robustly and systematically a model uses internal representations. For instance, a model simultaneously achieving perfect accuracy and high *conceptual instability* under small rephrasings or domain shifts indicates a shallow conceptual grounding. Conversely, a model with more modest performance but stable concept usage across transformations—like those described in §4—might be argued to exhibit deeper Computational Comprehension (Nanda et al., 2023).

From Output Scores to Conceptual Analysis

Traditional accuracy-based benchmarks measure *what* the model outputs but not *how* it arrives at those outputs. Mechanistic interpretability (§3) opens a window into *where* concept-like circuits reside (Olah et al., 2020), and *how* the model manipulates them. By integrating such insights into a new evaluation paradigm, researchers can design tasks and metrics that test for *invariance, compositional integration*, and *generalizable conceptual structure* rather than pure output correctness (Zahraei and Asgari, 2024). This reframes model assessment as probing the alignment between *internal computations* and *expected conceptual consistency*.

552

553

554

555

556

557

558

559

518

519

520

521

5.5 How Computational Comprehension Strengthens Model Evaluation

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

510

Focusing on Computational Comprehension offers a nuanced approach to evaluating model performance beyond traditional benchmarks. For instance, when Model A and B achieve identical scores, an analysis of Computational Comprehension may show that Model A employs identifiable concept circuits with robustness to synonymous phrases and logical coherence in reasoning, while Model B falters under minor perturbations and lacks clear internal structures, indicating that Model A demonstrates superior understanding despite similar benchmark scores. This approach reduces the risk of relying on superficial performance gains, as it requires models to demonstrate stable internal feature usage across a variety of tasks.

489 5.6 Addressing Objections and Feasibility

Critics of concept-centric metrics argue that imple-490 menting repeated transformations and deep abla-491 tion studies can be more resource-intensive than 492 numeric benchmarks, where a single accuracy or F1 493 score provides straightforward model comparison. 494 However, this added complexity is essential for ex-495 posing deeper strengths and weaknesses in model 496 cognition. Observational parallels can be drawn to 497 neuroscience and psychology, fields that similarly 498 adopt intricate testing to capture complex phenom-499 ena. Importantly, Computational Comprehension metrics do not displace existing evaluations; rather, 501 502 they complement them. Conventional benchmarks remain invaluable for quickly gauging baseline ca-503 pability or domain-specific performance. By layering in concept-invariance tests and interpretability measures, we can build a more complete picture of how robustly and structurally a model encodes the 507 knowledge it appears to exhibit. 508

6 Formal Framework for Assessing Computational Comprehension

511This section proposes a formal framework for as-512sessing Computational Comprehension, uniting513these ideas into a methodology that goes beyond514simple output metrics. In doing so, we shift em-515phasis from "How high is the accuracy?" to "How516stable and interpretable are the model's concept517representations?"

6.1 Concepts as Functional Mappings

We formalize a *concept* as a relation (or function) that maps a set of inputs \mathcal{X} to a corresponding set of outputs or labels \mathcal{Y} . Concretely, we write:

$$C\colon \mathcal{X} \to \mathcal{Y},$$
 522

where $C(x) \in \mathcal{Y}$ represents the ideal conceptual outcome (or property) for each input $x \in \mathcal{X}$. Crucially, C captures the *principle* or *logic* dictating how inputs relate to outputs, rather than a mere memorized pairing. For example, C(x) could designate whether x meets certain criteria, identifies a particular feature, or even transforms x according to a well-defined rule.

From this standpoint, a model M masters a concept C if it satisfies $M(x) \approx C(x)$ for all relevant inputs $x \in \mathcal{X}$. Such mastery demands more than just fitting a fixed set of training examples—it requires that M capture and robustly apply the underlying *principle* unifying all (x, C(x)) pairs. In other words, M should exhibit:

- **Consistency:** The model's predictions or outputs for new examples x' reflect the same structural or semantic rule encoded by C.
- Generality: The learned concept holds over the entire space \mathcal{X} (or its relevant subset), not solely for the specific instances seen during training.
- Stability: Small perturbations to x that preserve conceptual properties should not disrupt M's alignment with C, signaling that M latently encodes the concept rather than memorizing superficial patterns.

If the model genuinely *understands* a concept, it will consistently and accurately reproduce C(x)across transformations and contexts that preserve the conceptual core.

6.2 Transformational Robustness

One hallmark of true Computational Comprehension is the capacity to preserve conceptual validity under input transformations that leave the *concept* unchanged. Formally, suppose we define a set T of transformations $t: X \to X$ that maintain C(x) in the sense that:

$$C(t(x)) = C(x),$$
561

for all $t \in T$ and $x \in X$. If a model M genuinely grasps C, it should map t(x) to the same conceptual outcome as x—that is, $M(t(x)) \approx M(x)$ whenever C(t(x)) = C(x).

566

573

574

575

577

580

585

586

589

590

592

594

596

598

599

603

Benchmark Extension. Instead of scoring M solely on the raw input–output pairs (x, C(x)), we propose also testing across transformations (t(x), C(t(x))). The model's *robustness* score R(C, M) measures consistency under these transformations:

$$R(C,M) = \mathbb{E}_{(x,t) \in (X \times T)} \Big[\mathbf{1}(M(t(x)) = C(t(x))) \Big]$$

High R(C, M) suggests the model's internal representation is more than superficially memorizing (x, C(x)) pairs; it is capturing the general rule that remains invariant under t.

6.3 Linking Internal Representations via Mechanistic Interpretability

Even high robustness could, in principle, arise from sophisticated but purely *opaque* pattern-fitting. To address this risk, we invoke *mechanistic interpretability* (§3). Concretely, we hypothesize that for each concept C, there should exist a *sub-network* in M—some ensemble of neurons, attention heads, or circuits—that is critical to consistently mapping x to C(x). By identifying and testing that sub-network, we can verify whether:

- 1. Localization. Activations intensify in the subnetwork specifically when the model processes inputs related to *C*, compared to unrelated tasks.
- 2. Ablation Degrades C. Disabling or corrupting the targeted sub-network (e.g., *zero/mean ablation*) yields a sizable drop in R(C, M).
- 3. Consistency under Transformations. The same sub-network remains activated (or similarly structured) even for transformed inputs t(x), reinforcing that M is implementing C in a stable, generalizable manner.

These conditions connect conceptual invariance to an *observable* internal mechanism, offering a more holistic demonstration of Computational Comprehension.

604 6.4 Proposed Metrics and Procedures

We outline a possible experimental procedure, synthesizing ideas from previous sections: I. **Identify a Target Concept** *C***.** Formally define the function or relation you want to test (e.g., logical equivalence, color categorization, basic arithmetic).

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

- II. Generate Transformation Family T. Construct transformations $t \in T$ that preserve C, such as paraphrases, format changes, or domain shifts.
- III. Compute Baseline Accuracy. First assess M on pairs (x, C(x)) to ensure M handles straightforward versions of C acceptably.
- IV. Evaluate Robustness R(C, M). Arrange t(x) for $t \in T$ on test inputs, verifying if $M(t(x)) \approx C(t(x)) = C(x)$.

V. Mechanistic Analysis.

Find Relevant Circuits/Neurons. Use probing, circuit tracing, and feature visualization (see §3) to locate internal structures maintaining *C*.

Perform Ablation. Temporarily ablate (e.g. zero ablation) the identified sub-network. Does R(C, M) degrade substantially?

Test Activation Patterns. Confirm that these same sub-network activations systematically appear for transformed inputs t(x).

VI. Compare Models. If Models A and B both excel on raw benchmark scores, but Model A shows a more stable, localized concept subnetwork and robust R(C, M) under transformations, it can be argued to exhibit stronger Computational Comprehension of C.

6.5 Illustrative Example

Consider a concept $C_{\text{arithmetic}} = compute the sum of two numbers.$ Suppose our domain X is textual inputs like What is 123 + 456? We define T as transformations that (a) restate the query in different wording (Sum of 123 with 456), (b) embed the numbers in a story, or (c) add noise tokens that should not affect the sum. A model that *really* internalizes $C_{\text{arithmetic}}$ would consistently produce 579 even if the query is rephrased or placed in a distracting context.

Applying mechanistic interpretability, we might discover specific neurons or attention heads that reliably activate for arithmetic queries. Ablating them

699

should degrade addition performance across transformations. By empirically verifying such a subnetwork, we strengthen the claim that the model
holds a stable, concept-like mechanism for addition. Detailed comprehensive illustrations are in
Appendix C.

7 Conclusion

659

661

662

Have we truly advanced toward *machine understanding*, or merely built increasingly clever parrots? Can high scores on human-curated tasks ever capture the depth of a concept-driven architecture? And how do we know if a model's seeming mastery reflects stable conceptual knowledge rather than ephemeral statistical tricks? These questions anchor our reimagined path for AI research and practice, informed by a shift from anthropocentric definitions of understanding to more computationally grounded notions.

670 From Philosophical Inquiry to Applied671 Methodology

Throughout our exploration, we showed how philosophical debates regarding sensory modalities and 673 subjective experience caution us against overly gen-674 erous attributions of understanding. At the same 675 time, a strictly anthropocentric frame is unhelpful for analyzing computational systems unlike ourselves. By introducing Computational Comprehension, we direct focus to the operational consistency 679 of concepts represented in neural architectures-a definition testable in practice. Mechanistic inter-681 pretability then transforms what once seemed like philosophical abstraction into tangible questions about how models store, manipulate, and apply concepts.

Mechanistic Interpretability as Our Empirical Lens

688 Mechanistic interpretability serves as the bridge 689 between theoretical claims and observable model 690 behavior, exposing sub-networks or circuits most 691 relevant to a given concept. This approach lets us 692 analyze whether robust conceptual processing is 693 taking place, or whether we are witnessing super-694 ficial memorization patterns. When merged with 695 transformation tests, interpretability offers a more 696 holistic assessment of how well a model stabilizes 697 conceptual knowledge across contexts—an empiri-698 cal window into deeper AI cognition.

Reframing Benchmarks: Strengthening Our Evaluations

The community's current benchmark-focused environment, often driven by an arms race of difficulty, risks conflating performance with legitimate understanding once models surpass human scores. Instead, we propose measuring how well a model preserves conceptual invariants under transformations that should not alter conceptual meaning. This approach valuably shifts emphasis to *how* models achieve their results, rewarding internal coherence over rote data absorption.

Expanded Future Work in Multi-Modality

Our analysis focuses on large language models that primarily process text. However, multimodal systems integrating vision, audio, or other input streams introduce new challenges and opportunities for Computational Comprehension. Concepts rooted in visual or sensor-based data (e.g., object permanence, event segmentation) may be tested with transformations that alter viewpoints, backgrounds, or partial occlusions. At the same time, mechanistic interpretability becomes more complex, potentially demanding domainspecific circuit-mapping techniques. Leveraging these modalities could enrich how we define transformations and invariants, but the methodological burden of discovering and ablating relevant circuits also increases.

Taken together, these advances unveil a clearer vista for AI: a future where systems are judged not only by impressive performance metrics but by their demonstrated capacity to encode, transform, and reason with concepts in a sustained manner. By anchoring our evaluations in Computational Comprehension, complemented by mechanistic interpretability and transformation-based testing, we move beyond superficial achievements toward a deeper, more transparent, and ultimately more meaningful form of machine intelligence.

Limitations

Though this paper advocates for a more robust notion of Computational Comprehension and outlines an interpretability-based methodology for uncovering concept-centric representations, a number of limiting factors remain:

Scalability of Interpretability.Modern large745language models often feature hundreds of billions746

797

798

of parameters, making mechanistic analyses at neuron or circuit level extremely challenging. While
smaller models or targeted subsets of parameters
can be examined, interpreting a full-scale LLM remains computationally intensive and technically
difficult, potentially restricting such analyses to
well-resourced labs.

Concept Definition Dilemma. Despite the formal approach of modeling a concept as a function
or relation, actually specifying certain high-level or
context-dependent concepts (e.g., irony, advanced
logical reasoning, figurative language) is notoriously difficult. Oversimplified definitions risk overlooking emergent behaviors and interdependencies
that cannot be neatly captured by functional mappings alone.

Resource Constraints and Practical Feasibility.
Applying repeated transformations, running ablation experiments, and mapping concept circuits
can all be computationally expensive. This requirement may slow or hinder widespread adoption of
concept-based evaluations, particularly in industry
or smaller academic groups where hardware and
time are limited. While automated interpretability
pipelines are an active area of research, truly largescale application remains a formidable challenge.

Philosophical Caveats. Finally, even the most 773 thorough demonstration of stable concept circuits 774 does not prove the presence of conscious or sub-775 jective states. Our notion of Computational Com-776 prehension is deliberately non-anthropocentric; it 777 makes no claims about qualia or mental experi-778 ences. This helps avoid overattribution of human-779 like cognition to machines, yet it may leave some philosophical critiques unresolved, especially those 781 that posit experience as central to genuine understanding (e.g., Mary's Room).

> These limitations do not diminish the value of Computational Comprehension as a powerful conceptual and evaluative tool; rather, they highlight the need for further research, methodological refinement, and interdisciplinary collaboration in continuing to push the boundaries of *machine understanding*.

Ethics Statement

785

789

790

This work proposes a framework for examining and
measuring Computational Comprehension in large
language models using mechanistic interpretability
and transformation-based evaluations. While our

aim is to foster deeper insights into a model's conceptual capacities, we acknowledge several ethical considerations:

- Misuse of Interpretability. Enhanced methods for dissecting internal model mechanisms, though beneficial for transparency, may be misapplied to extract proprietary information or target sensitive aspects of model behavior. Researchers and developers should exercise care in deciding which details are publicly disclosed, balancing transparency with the potential for malicious exploitation.
- **Bias and Fairness.** As we refine tools to identify conceptual circuits, it is equally important to ensure they detect and mitigate biases. Models can unintentionally encapsulate harmful stereotypes that escape surface-level benchmarking. Mechanistic interpretability thus has an ethical imperative to uncover and address such issues more systematically.
- Data Privacy. Many interpretability protocols rely on intensive probing or ablation experiments with curated datasets. Researchers should remain cognizant of privacy considerations and legal constraints—particularly when dealing with personal or sensitive data—and follow best practices to anonymize or redact identifying information.
- Anthropomorphizing and Responsibility. While Computational Comprehension provides a non-anthropocentric lens, there remains a risk of over-attributing agency or accountability to models. We clarify that these evaluations do not confer moral or legal responsibility on the system and should not be used to justify autonomy or moral standing for LLMs.
- **Potential for Unequal Access.** Robust interpretability frameworks can be computationally and logistically expensive, raising concerns that only a small segment of the AI community will have sufficient resources to conduct these studies. The field should encourage shared resources, open-source tools, and collaboration to ensure that concept-based evaluations do not reinforce inequitable access or hinder responsible innovation.

By presenting Computational Comprehension as an additional source of insight, our goal is to

complement, rather than displace, existing ethical and regulatory frameworks. We advocate for open, interdisciplinary dialogue on how best to leverage detailed model understanding to enhance both performance and societal benefit, mindful of the risks of misuse, bias, or inequity. 850

References

851

853

861

874

875

876

877

885

890

896

- Daman Arora, Himanshu Gaurav Singh, et al. 2023. Have llms advanced enough? a challenging problem solving benchmark for large language models. arXiv preprint arXiv:2305.15074.
 - Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big? FAccT.
- Michael Timothy Bennett. 2023. On the computation of meaning, language models and incomprehensible horrors. In International Conference on Artificial General Intelligence, pages 32-41. Springer.
- Carl Bereiter. 2002. Education and Mind in the Knowledge Age. Routledge.
- Leonard Bereska and Efstratios Gavves. 2024. Mechanistic interpretability for ai safety - a review. arXiv.
- Nigel Collier, Fangyu Liu, and Ehsan Shareghi. 2022. On reality and the limits of language data: Aligning llms with human norms. In Annual Meeting of the Cognitive Science Society.
- Lucia Foglia and Robert Wilson. 2013. Embodied cognition. Wiley Interdisciplinary Reviews: Cognitive Science, 4.
- Clémentine Fourrier, Nathan Habib, Alina Lozovskaya, Konrad Szafer, and Thomas Wolf. 2024. Open llm leaderboard v2. https://huggingface. co/spaces/open-llm-leaderboard/open_llm_ leaderboard.
- Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever, Jan Leike, and Jeffrey Wu. 2024. Scaling and evaluating sparse autoencoders. arXiv preprint arXiv:2406.04093.
- Sachin Goyal, Ziwei Ji, Ankit Singh Rawat, Aditya Krishna Menon, Sanjiv Kumar, and Vaishnavh Nagarajan. 2023. Think before you speak: Training language models with pause tokens. arXiv preprint arXiv:2310.02226.
- Stefan Heimersheim and Neel Nanda. 2024. How to use and interpret activation patching. arXiv preprint arXiv:2404.15255.
- Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2020. Measuring massive multitask language understanding. arXiv preprint arXiv:2009.03300.

Jos ´e Hern ´andez-Orallo. 2020. Ai evaluation: On bro-	897
ken yardsticks and measurement scales.	898
Frank Jackson. 1986. What mary didn't know. <i>The Journal of Philosophy</i> .	899 900
Kevin Meng, David Bau, Alex Andonian, and Yonatan	901
Belinkov. 2022. Locating and editing factual associ-	902
ations in gpt. <i>Advances in Neural Information Pro-</i>	903
<i>cessing Systems</i> , 35:17359–17372.	904
Raphaël Millière and Charles Rathkopf. 2024. Anthro-	905
pocentric bias and the possibility of artificial cogni-	906
tion. In <i>ICML 2024 Workshop on LLMs and Cogni-</i>	907
<i>tion</i> .	908
Philipp Mondorf and Barbara Plank. 2024. Beyond	909
accuracy: Evaluating the reasoning behavior of large	910
language models - a survey. <i>ArXiv</i> , abs/2404.01869.	911
Neel Nanda, Lawrence Chan, Tom Lieberum, Jesse	912
Smith, and Jacob Steinhardt. 2023. Progress mea-	913
sures for grokking via mechanistic interpretability.	914
<i>ICLR</i> .	915
Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel	916
Goh, Michael Petrov, and Shan Carter. 2020. Zoom	917
in: An introduction to circuits. <i>Distill</i> .	918
Alethea Power, Yuri Burda, Harri Edwards, Igor	919
Babuschkin, and Vedant Misra. 2022. Grokking:	920
Generalization beyond overfitting on small algorith-	921
mic datasets. <i>arXiv preprint arXiv:2201.02177</i> .	922
P Rajpurkar. 2016. Squad: 100,000+ questions for machine comprehension of text. <i>arXiv preprint arXiv:1606.05250</i> .	923 924 925
John R Searle. 1980. Minds, brains, and programs.	926
Behavioral and Brain Sciences.	927
Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao,	928
Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch,	929
Adam R Brown, Adam Santoro, Aditya Gupta,	930
Adrià Garriga-Alonso, et al. 2022. Beyond the	931
imitation game: Quantifying and extrapolating the	932
capabilities of language models. <i>arXiv preprint</i>	933
<i>arXiv:2206.04615</i> .	934
James Strachan, Dalila Albergo, Giulia Borghini, Oriana	935
Pansardi, Eugenio Scaliti, Saurabh Gupta, Krati Sax-	936
ena, Alessandro Rufo, Stefano Panzeri, Guido Manzi,	937
Michael Graziano, and Cristina Becchio. 2024. Test-	938
ing theory of mind in large language models and	939
humans. <i>Nature Human Behaviour</i> , 8:1–11.	940
 Adly Templeton, Tom Conerly, Jonathan Marcus, Jack	941
Lindsey, Trenton Bricken, Brian Chen, Adam Pearce,	942
Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy	943
Cunningham, Nicholas L Turner, Callum McDougall,	944
Monte MacDiarmid, C. Daniel Freeman, Theodore R.	945
Sumers, Edward Rees, Joshua Batson, Adam Jermyn,	946
Shan Carter, Chris Olah, and Tom Henighan. 2024.	947
Scaling monosemanticity: Extracting interpretable	948
features from claude 3 sonnet. <i>Transformer Circuits</i>	949
Thread.	950

Alex Wang. 2018. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*.

951

952

957

960

961

962

964

965

966

967

970

971

973

974

975

977

978

979

981

991

993

997

998

1001

- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019. Superglue: A stickier benchmark for general-purpose language understanding systems. *Advances in neural information processing systems*, 32.
- Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*.
- Zeping Yu and Sophia Ananiadou. 2024a. Interpreting arithmetic mechanism in large language models through comparative neuron analysis. *arXiv preprint arXiv:2409.14144*.
 - Zeping Yu and Sophia Ananiadou. 2024b. Neuronlevel knowledge attribution in large language models. In Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing, pages 3267–3280.
- Pardis Sadat Zahraei and Ehsaneddin Asgari. 2024. Turingq: Benchmarking ai comprehension in theory of computation. *EMNLP Findings*.
- Wei Zhang, Chaoqun Wan, Yonggang Zhang, Yiu-ming Cheung, Xinmei Tian, Xu Shen, and Jieping Ye. 2024. Interpreting and improving large language models in arithmetic calculation. arXiv preprint arXiv:2409.01659.
- Kun Zhou, Yutao Zhu, Zhipeng Chen, Wentong Chen, Wayne Xin Zhao, Xu Chen, Yankai Lin, Ji-Rong Wen, and Jiawei Han. 2023. Don't make your llm an evaluation benchmark cheater. *arXiv preprint arXiv:2311.01964*.

Appendix

A Preliminary Exploration of Transformations for Logic

For logical inference tasks, we define a relation $C(\phi, \psi)$ as a function returning true if ψ is logically entailed by ϕ in propositional logic. To test transformational robustness, one can apply rewritings or reorder terms without altering the underlying logical dependencies. Examples include variable renaming, distribution of negations, or rewriting implications (e.g., $p \rightarrow q$ as $\neg p \lor q$). If M displays high consistency under such transformations, it suggests the model holds a more generalizable representation of logical entailment, instead of merely memorizing surface patterns.

B Further Discussion of Mechanistic Ablation Techniques

Main text discussions of ablation (Section 3) can be1004expanded by referencing specialized protocols such1005as Neuron-Level Knowledge Attribution in Large1006Language Models. These involve:1007

1003

1014

1015

1016

1017

1018

1019

1020

1021

1022

1023

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1039

1040

1041

1042

- 1. Identifying neuron groups strongly correlated1008with a target concept based on their activation1009patterns.1010
- 2. Zeroing out or injecting noise into those neuron activations while the remaining network remains intact.10111012
- 3. Measuring the consequent drop in performance on tasks or test suites that rely on the targeted concept.

Marked degradation suggests that knowledge of the concept is at least partially localized in those neurons. However, the possibility of distributed or redundant representations must be carefully weighed, potentially requiring multiple ablation variants or iterative neuron-discovery steps.

C Illustrative Example: Arithmetic Comprehension

Consider a concept $C_{\text{arithmetic}} = compute the sum of two numbers.$ This concept provides an ideal test case for Computational Comprehension due to its well-defined structure and clear expected outputs.

Concept Definition and Input Domain

Our domain X consists of textual inputs requesting arithmetic operations, such as "What is 123 + 456?". The concept function $C_{\text{arithmetic}}$ maps these inputs to their numerical solutions (in this case, 579).

Transformation Family

We define transformation family T to include operations that preserve the underlying arithmetic concept while varying surface presentation:

- 1. **Paraphrasing:** Restate the query using different wording ("*Calculate the sum of 123 and 456*", "What do you get when you add 123 to 456?")
- 2. Contextual Embedding: Place the arithmetic operation within a narrative context ("John had 123 apples and received 456 more. How many apples does John have now?")
 1043

- 3. Noise Addition: Insert irrelevant tokens that should not affect the computation ("What is, um, let me think, 123 plus, you know, 456?")
 - 4. Format Variation: Present numbers in different formats ("What is one hundred twentythree plus four hundred fifty-six?")

Evaluating Computational Comprehension

A model that truly comprehends the concept of addition would consistently produce the correct sum (579) across all these transformations. This consistency indicates that the model has internalized the abstract operation of addition rather than merely memorizing specific input-output patterns.

Mechanistic Analysis Protocol

1047

1048

1049

1050

1051 1052

1053

1054

1055

1056

1058

1059

1060

1061

1062

1063

1064

1066

1067

1068

1069

1071

1073

1074

1075

1076

1077

1078

1079

1080

1082

1083

1084

1085

1086

1087 1088

1089

1090

1092

Recent work in mechanistic interpretability has demonstrated methods for identifying arithmetic circuits in language models. We propose the following protocol to verify computational comprehension of arithmetic:

- 1. **Circuit Identification:** Using causal mediation analysis or activation patching, identify the specific neurons, attention heads, or MLP modules that activate when processing arithmetic queries. Research suggests these components often reside in mid-sequence early layers and the final token's later layers.
- 2. Activation Pattern Analysis: Record activation patterns across the transformation family T. A model with true computational comprehension should show consistent activation in the same sub-networks regardless of how the arithmetic query is presented.
- 3. **Targeted Ablation:** Temporarily disable the identified arithmetic circuit components through zero ablation or mean ablation. If these components are truly responsible for arithmetic processing, performance should degrade specifically on arithmetic tasks while leaving other capabilities intact.
- 4. Cross-Task Comparison: Compare activation patterns during arithmetic tasks with patterns during number retrieval tasks and factual knowledge questions. True arithmetic comprehension should exhibit distinct circuit activation compared to mere number recognition or general knowledge retrieval.

Expected Observations	1093
If a model possesses computational comprehension	1094
of arithmetic, we would expect to observe:	1095
1. Consistent performance $(M(t(x)) \approx C(x) =$	1096
579) across all transformations $t \in T$	1097
2. Identifiable sub-networks that activate specifi-	1098
cally for arithmetic operations	1099
3. Significant performance degradation on arith-	1100
metic tasks (but not other tasks) when these	1101
sub-networks are ablated	1102
4. Similar activation patterns across different pre-	1103
sentations of the same arithmetic problem	1104
This arithmetic example demonstrates how Com-	1105
putational Comprehension goes beyond surface-	1106
level accuracy to examine the internal mechanisms	1107
supporting concept representation and manipula-	1108
tion. By verifying both transformational robustness	1109

tion. By verifying both transformational robustness1109and the presence of dedicated conceptual circuits,1110we can distinguish models that truly comprehend1111arithmetic from those that merely pattern-match on1112specific input formats.1113