. janicre: A Log-Scale, Reversible Semantic-Commit Manifest for Multi-Agent Software Reasoning

Anonymous ACL submission

Abstract

002 Large language models (LLMs) promise repository-scale assistance, yet real projects routinely exceed 128 k-token windows and fragment key logic across heterogeneous files. Existing techniques-RAG pipelines, chunked retrieval, or brute-force long-context prompting-either leak code, stumble on cross-file reasoning, or incur quadratic cost.

> We realise . janicre via an arbitrary-depth abstraction pipeline (Stage $1 \dots k$), whose cumulative JSON snapshots plus an HTML semantic map form the final manifest; stopping at any depth trades tokens for detail while the worst-case growth remains $\Theta(\log N)$.

Beyond human-to-LLM use, .janicre functions as a model-agnostic exchange layer: distinct agents (GPT-40, Claude, Gemini, etc.) can inspect the same manifest, attach thinking_trace logs, and negotiate edits-enabling true agent-to-agent (A2A) collaboration without exposing raw source. Coupled with git-like semantic diffs, the manifest becomes a living memory of design rationale that is auditable by both machines and humans.

015

017

018

022

027

We formalise the schema, analyse its compression bounds, and outline an empirical plan—HumanEval, MBPP. and delta-stability benchmarks-to compare .janicre-augmented reasoning against standard retrieval and full-file prompting.

. janicre thus upgrades IR-level code summarisation into a universal, dialogueready substrate for large-scale, multi-agent software intelligence.

.janicre abstracts multi-**Highlight:** language codebases into a $\log N$ -scale manifest that aligns with transformer reasoning and agent 039 collaboration-without code leakage. It bridges human design intent and LLM reasoning via a real-time semantic interface. 042

Introduction 1

Problem Context 1.1

LLM assistants have two systemic bottlenecks: (i) token budgets (>128 k tokens are still dwarfed by modern repositories) and (ii) fragmentationdomain logic is scattered across files, packages, and services. Current practice therefore relies on ad-hoc snippets or retrieval pipelines, risking both context loss and code leakage.

1.2 Compression Inspiration: From MP3 and MP4 to Code Semantics

Just as MP3 compresses audio by removing frequencies imperceptible to humans, and MP4 compresses video by encoding delta frames and keyframes, .janicre compresses code by extracting only semantically relevant units. These units represent the purpose and behavior of the system in a compact and interpretable form. This mirrors the design goal of enabling LLMs to focus on core logic rather than superficial syntax, providing an efficient input representation aligned with human intuition and transformer attention.

1.3 Limitations of Existing Approaches

Retrieval-augmented generation (RAG) retrieves relevant snippets (Lewis et al., 2020) but struggles with cross-file reasoning due to limitations in context aggregation, such as chunk-based retrieval that fragments inter-file dependencies. Longcontext models (Dong et al., 2023) reduce manual selection but incur quadratic computational costs in transformer attention and degrade on gigabytescale repositories due to inefficiencies in scaling long sequences.

1.4 Transformer-centric Abstraction

Rather than a fixed three-step recipe, we allow a variable-depth cascade: Stage 0 (raw code) \rightarrow Stage 1 (AST) \rightarrow Stage 2 (purpose) $\rightarrow \cdots \rightarrow$

043

044

045

- 068 069 070 071
- 072

073

074

075

076

078

079



Figure 1: High-level overview of .janicre: (left) multi-stage semantic compression with $\Theta(\log N)$ attention; (right) mapping to the classical SDLC V-model.

Stage k (intent, theory). Each step retains only the *attention-anchoring tokens* (names, types, keywords), analogously to how MP3 preserves auditory maskers. Each stage is stored as a reversible semantic commit, enabling loss-bounded roundtrip across abstraction depths. Because k can be capped at $O(\log N)$, the pipeline still yields $\Theta(\log N)$ tokens in the worst case, yet lets practitioners exit early when smaller windows suffice.

To enhance transformer alignment, we define a soft attention prior between semantic units based on their abstraction depths. Specifically, we compute the attention weight between unit i and j as:

$$w_{i,j} = \exp(-\alpha \cdot |d_i - d_j|), \quad A_{i,j} = \frac{w_{i,j}}{\sum_j w_{i,j}}$$

where d_i denotes the abstraction depth of unit *i*, and α is a tunable decay factor. This encourages intra-stage attention and smooths cross-stage transitions, aligning the attention pattern with the underlying semantic hierarchy.

1.5 Our Contribution

090

094

100

101

A Unified Time-Aware Interface for LLM Reasoning

Unlike static representations, .janicre introduces a time-aware, dialogue-based manifest that
models not just software structure, but its design rationale across time. This enables LLMs
to interpret codebases as evolving conversations—aligning each unit with a reasoning turn and associated expectation.

107

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

Moreover, .janicre supports *dynamic compression control*, allowing the manifest to adapt its granularity based on context: high-detail traces for chain-of-thought reasoning, and compact structural summaries for fast inference or retrieval. This flow-sensitive abstraction turns .janicre into a scalable LLM interface, bridging software design, memory, and attention cost.

Together, these features establish .janicre as a temporal-semantic interface between software systems and large language models.

1.6 Our Contribution

- 1. Schema. We define .janicre, a manifest that compresses repository intent into $\mathcal{O}(\log N)$ tokens.
- 2. **CLI.** An open-source generator produces manifests for Python / C++ projects; Rust/Go plugins are planned.
- 3. **Security.** Only derived metadata crosses the firewall, preventing proprietary code exposure.

2 Prompt-oriented DSLs

Syntactic IRs. Early LLM-for-code studies rely on *syntactic* representations such as ASTs (Code-BERT) or CFGs (GraphCodeBERT). These IRs expose parse-tree structure but seldom encode *intent, semantics*, or cross-module relationships.

Prompt-oriented DSLs. More recent work shifts to prompt engineering DSLs (e.g., agent scripts, Toolformer graphs). While flexible, these DSLs are *imperative* and often tied to a specific runtime, limiting reuse across projects.

. janicre: a semantic, declarative IR. Our schema occupies the gap between the two trends:

- Semantic layer captures *purpose*, *I/O contracts*, and *module hierarchy*, not just syntax.
- Declarative form JSON/YAML manifest aligns with transformer attention, enabling token-efficient reasoning ($O(\log N)$ tokens).
- LLM-friendly graph unit-level nodes and explicit edges mirror token-level attention (Vaswani et al., 2017).

Recent frameworks such as PDL (Chen et al., 2024), LangGPT (Wang et al., 2024), and Doc-CGen (Pimparkhede et al., 2024) define declarative layers for modular LLM interaction and document-conditioned code generation. These approaches improve prompt composability, yet typically remain imperative, runtime-tied, or file-local. In contrast, semantic-commit and timeaware intermediate representations—structured to persist abstraction history—remain largely unexplored.

The .janicre schema is not merely a descriptive manifest, but a structurally optimized intermediate format for LLM-based reasoning.

First, we choose JSON/YAML as the base format for the following reasons:

- LLM Compatibility: Transformer-based LLMs attend over token sequences, and key-value structured formats (like JSON) yield stable and predictable attention paths, helping models form logical groupings during inference.
- Multilingual Metadata Extraction: Languages such as Python, JavaScript, Rust, or Go expose ASTs that can be naturally mapped to JSON, enabling unified structural representation.
- Declarative Abstraction: Unlike imperative code, this schema allows systems to be described independently of runtime execution—capturing semantics and modular relationships.

Moreover, we treat all software as ultimately a composition of input-output transformations. This aligns with the principle that LLMs operate effectively when provided with structured descriptions of units, each explicitly named and ordered. For example:

"units": ["startGame", "endGame", "
spawnEnemy"]
''' A simplified excerpt is shown
below:
```json
"units": ["startGame", "endGame", "
spawnEnemy"],
"commit_meta": {
"timestamp": "2025-05-20T12:34Z",
"depth": 3
}

Such structure enables LLMs to attend across unit boundaries and predict interactions more accurately.

This design echoes attention-based reasoning, as introduced in (Vaswani et al., 2017), and adapts it to LLM–software interfaces. Rather than prompt full source files, we create interpretable hooks for the model to follow.

Moreover, while the core .janicre schema is formally defined, we consider any simplified variants (e.g., flat representations), partial modular forms, or implementation-specific adaptations to be legitimate extensions under the same conceptual framework. These derived or restructured formats remain within the scope of our proposal as long as they preserve the semantic objectives and hierarchical abstraction principles defined herein.

### 2.1 Preliminary Observation

A manually structured Shogi engine (²GB) suggests a **97% token reduction** when encoded via .janicre, compared with full-file prompting.

Schema Overview. The .janicre manifest describes a software system's purpose, components, and logic in a JSON/YAML format optimized for LLM reasoning. Each manifest includes fields such as global_objective, code_sharing, and code_logic, which summarize the architecture, I/O contract, and computational units.

A simplified excerpt is shown below:

"units":	["startGame",	"endGame",	"
spawnEnemy"]			

The manifest's complete structure—including global_objective, code_sharing, and

detailed code_logic units—is provided in the
supplementary file example.janicre.json.

238

239

241

242

243

246

247

248

251

257

258

261

263 264

265

269

271

272

273

274

277

278

281

Modular Composition. The .janicre schema is inherently composable. Multiple sub-manifests—e.g., ui.janicre for HTML layout and backend.janicre for Python/JS logic—can be referenced or included by a higherlevel core.janicre. This enables distributed teams or tools to independently define interface, logic, and algorithms, yet produce a unified LLM-interpretable manifest. A merging protocol resolves shared objectives and unifies component trees and unit indexes across manifests.

> Merge respects commit order; higher-depth units are deterministically regenerated from lowerdepth blobs, guaranteeing reversible abstraction quality.

#### 2.2 From Git Commits to Dialogue Units

Traditional version control systems like Git decompose software into a series of commits—each capturing a structural diff and accompanied by a descriptive message. These commits form a linear or branching history, enabling developers to track the evolution of code over time.

. janicre draws inspiration from this idea, but shifts the focus from syntactic diffs to semantic intent. Each unit in a manifest acts like a commit: it is named, scoped, and justified. The optional thinking_trace serves as a semantic commit message, recording the rationale, prompt context, and decision steps behind the unit's creation or modification.

Unlike Git, which records lines of code, .janicre records units of purpose and reasoning. This structure provides an audit trail for LLMs—not just for what changed, but why—and enables model-aware tools to interpret, revise, or simulate development at the semantic level.

In this sense, .janicre is a Git-like protocol for software semantics, optimized for transformerbased reasoning systems. It captures code not merely as syntax but as a living dialogue of design decisions.

### 2.3 Conversational Interface Inspiration

Why Meaning Emerges Only in Dialogue. Transformer-based LLMs operate by predicting the next token, but prediction alone does not guarantee semantic coherence. Given an underspecified prompt such as "Write about time", RLHF- tuned models tend to regress toward a *single*, *safest average* response that maximizes annotator agreement. This phenomenon—semantic convergence toward a median output—arises precisely because the input lacks *directional intent*. The less a prompt specifies expectation, the more the model collapses meaning into neutrality.

285

287

290

291

293

294

295

296

297

299

301

302

303

304

305

306

307

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

332

Dialogue format resolves this degeneracy. Each conversational turn embeds an expectation shaped by context: the prior utterance narrows the manifold of plausible continuations and forces the model to align with a specific semantic target. **Only in dialogue can an LLM both imagine and justify a precise value—because dialogue provides expectation.** We thus view conversation as a *convergence protocol*—a structure that transforms token-level prediction into context-aligned reasoning.

**Structural Properties of Dialogue.** Beyond convergence, dialogue also supplies inductive biases beneficial to transformer models:

- **Turn-taking** constrains the scope per utterance, preventing attention from overextending.
- **Intent segmentation** frames each utterance as a discrete, nameable semantic unit.
- **Contextual grounding** maintains a shared, evolving state across steps.

These properties are not cosmetic—they *scaf-fold* model behavior in alignment, memory, and structured inference.

**Empirical Support.** This perspective is supported by empirical studies. Ouyang et al. demonstrate that instruction-tuning via conversational prompts improves model stability and alignment with human preferences (Ouyang et al., 2022). Christiano et al. show that turn-based feedback via reinforcement learning yields more robust policies than one-shot or batch supervision (Christiano et al., 2017). These results underscore that conversation is not merely interaction—it is optimization.

**Design Implication for .janicre.** This insight motivates the structural choices in .janicre. Each unit—explicitly named, scoped by input/output, and embedded in a hierarchical context—acts as a semantic utterance in a system-level dialogue. Rather than presenting code as flat syntax, .janicre frames it as

extensions may support Rust, Go, and declarative

384

386

389

390

391

392

393

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

intent-driven discourse. Software becomes a dialogue among components, and LLMs can reason over it not as text, but as communicative structure.

333

334

335

337

339

340

341

345

347

350

351

357

361

363

364

367

371

374

Why Meaning Emerges Only in Dialogue. Transformer-based LLMs operate by predicting tokens—but prediction alone does not guarantee semantic coherence. In unconstrained text, the output space is vast and amorphous; without external structure, the model's generation lacks a center of gravity.

Conversational format resolves this. Each turn imposes an expected response, compressing the space of plausible outputs and forcing the model to converge toward semantically meaningful values.
It is only within dialogue that an LLM can both imagine and justify a specific value—because dialogue embeds expectation.

In essence, the act of dialogue is a convergence protocol: a structural constraint that transforms token-level prediction into intent-aligned imagination. We posit that *conversation is not merely a UX design, but the minimal viable structure under which large language models can produce grounded meaning.* 

This motivates our design of .janicre. By treating each software unit as a scoped utterance—with name, purpose, and contract—it mirrors the same structure that enables LLMs to think. Software becomes a dialogue between components, and LLMs can reason about it not as code, but as communicative structure.

# 3 Reference Implementation & Multi-Stage Pipeline

# 3.1 Pipeline Overview

We implement a generator janicre-cli that emits *Stage-i* JSON files (i = 1...k), each representing a progressively abstracted view of the codebase—from raw ASTs to high-level semantic units. This cascade allows developers to tune token granularity, stopping at any Stage-*j* to fit context limits while preserving structural clarity.

375Language-Agnostic Design.Unlike language-376specific tools, janicre is not bound to any par-377ticular syntax or ecosystem.It supports Python,378JavaScript, HTML, Vue, and even domain-specific379DSLs.Each Stage-i manifest shares a unified380schema capturing purpose, input/output contracts,381and modular hierarchy, enabling LLMs to reason

**HTML Semantic Map.** All Stage-*i* JSONs are merged into an interactive HTML canvas powered by Vue. Clicking a unit reveals its source_ref, editing its purpose live-updates the underlying JSON, and exporting yields a consolidated .janicre. This interface offers both machineparsable structure and human-navigable semantics—bridging reasoning and authoring within a single workflow.

consistently across heterogeneous systems. Future

formats like YAML or Terraform.

# 4 Semantic Commits & Thought Trace

To extend .janicre beyond static manifest, we introduce thinking_trace—a per-unit log of reasoning steps and LLM prompts that shaped each commit.

This allows:

- traceable origin of unit definitions (e.g., spawn logic);
- comparison across models (e.g., GPT vs Claude);
- intent continuity under refactoring or future evolution.

Just as Git preserves structural diffs, .janicre maintains conceptual diffs across units. This enables semantic scaffolding for LLM pretraining, where prompt-response pairs can serve as alignment records.

**Git-style Correspondence.** We model each unit as a semantic commit, allowing .janicre to function as a time-aware versioning interface for LLMs. The correspondence is as follows:

- Commit message  $\rightarrow$  thinking_trace rationale
- Diff  $\rightarrow$  semantic delta between abstraction depths
- Timestamp  $\rightarrow$  temporal ordering of abstractions

By merging structure with discourse,421.janicre becomes not just a snapshot—but a422living memory of design evolution.423

492

493

494

495

496

497

498

499

502

503

504

505

506

508

510

511

512

513

514

515

516

517

518

470

471

# 5 Token and Security Analysis

424

425

426

427

428

429

430

431

432

433

434

435

436

437

438

439

440

441

442

443

444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

Let c (0 < c < 1) be the empirical geometric decay factor, empirically observed as  $c \approx 0.25$  in our corpus. We model the token count at abstraction depth *i* as:

$$T_i = c^i N$$

where N is the original token count at Stage 0. This geometric decay reflects progressive semantic pruning. At the maximum depth k, the token count converges to:

$$T_k \approx \Theta(\log N), \text{ for } k = \left\lceil \log_{1/c}(N) \right\rceil$$

Practitioners can choose any  $j \leq k$  depending on context window constraints. Given a maximum token budget B (e.g., 32k tokens for GPT-4), the minimum required stage depth  $k^*$  can be computed as:

$$k^{\star} = \left\lceil \log_{1/c} \left( \frac{N}{B} \right) \right\rceil$$

This enables automated tradeoffs between semantic fidelity and budget, making the system adaptive to various model capacities.

Notably, the compression is not merely a pruning—it preserves transformer-aligned units such as attention anchors and semantic landmarks, retaining the system's functional semantics at every level.

# 6 Planned Evaluation

We will benchmark HumanEval and MBPP to measure (i) accuracy gain versus snippet baselines, (ii) token savings, and (iii) zero-leakage assurance by manual review.

Strategic Value of Model-Agnostic Benchmarking. By aligning .janicre evaluation with standard LLM benchmarks such as SWE-bench and ToolBench, we enable consistent, modelagnostic assessment of code reasoning performance. Crucially, this compatibility allows users to switch between LLM backends (e.g., GPT-40, Claude, Gemini) without restructuring the prompt interface or modifying the underlying manifest.

This interchangeability reduces vendor lock-in and amortizes the cost of adopting .janicre—transforming initial investment in manifest structuring into a long-term asset reusable across evolving model ecosystems. We argue that this flexibility plays a pivotal role in enterprise-scale LLM deployment, where rapid model turnover and API discontinuities are frequent. As such, benchmark alignment thus becomes not merely an evaluation methodology, but a cornerstone of long-term operability, cost mitigation, and architectural agility in evolving LLM ecosystems.

**Delta Stability under Code Evolution.** Inspired by MP4's delta frame model, we plan to measure the stability of .janicre manifests across codebase changes. Specifically, we will benchmark whether small code edits (e.g., bug fixes or new functions) result in local diffs within the manifest. This helps evaluate the long-term viability of .janicre as a persistent abstraction layer over evolving repositories.

We will also benchmark abstraction round-trip fidelity (shallow  $\rightarrow$  deep  $\rightarrow$  shallow) to quantify reversible commit accuracy and ensure semantic equivalence across depth transformations.

**Depth-Sweep Study.** For  $j = 1 \dots k$ , we feed {Stage  $\leq j$ } into GPT-40-32k and measure task accuracy against token count. This allows us to verify empirical geometric decay, validate compression efficiency, and identify the "elbow" depth—i.e., the optimal tradeoff point between token budget and semantic fidelity.

### Limitations

While . janicre enables efficient representation and reasoning, several limitations remain.

First, it is designed for static manifest generation. Dynamic behaviors—such as eval-based execution, reflection via getattr, or runtime mutation of control flow—cannot be fully captured through static analysis alone. As such, runtime-specific logic may require complementary trace-based instrumentation.

Second, current support is limited to statically typed or semi-structured languages; dynamically reflective systems may require additional adaptation for full semantic fidelity.

Third, the compression quality depends on developer-authored metadata such as docstrings and function modularity. Poorly documented or monolithic code may yield underspecified or noisy manifests, reducing the effectiveness of downstream LLM reasoning.

Finally, .janicre guarantees loss-bounded round-trip across depths, though regenerating full source from the highest abstraction may be lossy

520

521

522

524

525

527

528

530

531

532

533

535

536

539

540

541 542

544

547

548

549

552

554

557

560

561

564

in syntactic details.

# 7 Potential Risks

We discuss three main risks. (1) Informationleakage risk: even though .janicre removes raw source, a manifest may still expose sensitive architectural metadata; we propose configurable redaction levels and cryptographic hashes to mitigate this. (2) Misuse by malicious agents: an attacker could inspect the manifest to craft targeted exploits; we suggest access-control wrappers when .janicre is generated from proprietary code. (3) Over-reliance on compressed views: developers might miss security-critical edge-cases that were pruned; we therefore recommend fall-back to deeper stages or full static analysis for safety-critical code.

# 8 Conclusion

. janicre compresses software intent and structure into a single manifest, enabling LLMs to assist on large, fragmented codebases within token and privacy constraints.

Beyond practical usage, the schema may serve as a candidate format for future LLM pretraining corpora. With sufficient .janicre samples accumulated across projects, models could learn modular abstraction, intent classification, and high-level architectural reasoning—analogous to how code search engines or doc2vec embeddings are built from structured repositories today.

We are drafting a formal JSON Schema for .janicre v2.1 to enable validation and integration into IDE tooling.

Additionally, by abstracting software logic into a model-independent format and aligning it with benchmark-based evaluation, .janicre ensures long-term interoperability across evolving LLM architectures. This enables seamless model substitution, minimizes prompt engineering overhead, and mitigates risks of vendor lock-in or platformspecific obsolescence.

In doing so, .janicre delivers not only structural and reasoning efficiency, but also economic sustainability—transforming initial LLM integration efforts into a durable, reusable asset that preserves value amid a rapidly shifting model landscape.

#### Appendix A: . janicre Manifest 565 **Blueprint** 566 The fragment below illustrates the full conceptual 567 scaffold that .janicre is designed to capture. 568 It can be read as a "check-list" for repository-to-569 LLM interaction. 570 global_objective • **purpose**: compute closed-form solutions to 573 symmetrical tiling problems 574 • **technical**_category discrete 575 mathalgorithms, CLI scripts 576 • code_sharing • language: language-agnostic (supports Python, 579 JavaScript, HTML, etc.) 580 structure.external_links : none 581 • code_logic overview.core_algorithm 584 ٠

combinatorial enumeration under inversion symmetry overvie modular decomposition with separation of UI and logic 586

Frontend:	587
• module: ui.vue	588
• component: gameUI	589
• units: startGame, endGame, spawnEnemy	590
Unit Definitions:	591
• name: startGame	592
– type: method	593
- purpose: Initialize and start game loop	594
– detail.input: -	595
<ul> <li>detail.output: Running game state</li> </ul>	596
– detail.steps:	597
1. Initialize game flags and score	598
2. Reset player and bullet arrays	599
3. Setup and play background music	600
4. Launch game loop with	601
requestAnimationFrame	602
– detail.formula: -	603
– detail.theory: UI state initialization and ani-	604
mation loop	605
– detail.application: Browser-based arcade	606
games	607
• name: endGame	608
– <b>type</b> : method	609
<ul> <li>purpose: Stop the game and display results</li> </ul>	610
– detail.input: -	611

612	- detail.output: Stopped state and final score	Sa
613	shown	
614	– detail.steps:	
615	1. Pause background music	
616	2. Set gameOver = true	
617	3. Show result UI	A
618	– detail.formula: -	
619	- detail.theory: State transition in game UI	
620	- detail.application: Game-over screen logic	
621	• <b>name</b> : spawnEnemy	M
622	– type: method	141
623	- purpose: Create a new enemy on canvas	
624	– detail.input: -	
625	- detail.output: Enemy rendered in game	
626	space	
627	– detail.steps:	
628	1. Randomize spawn coordinates	
629	2. Push enemy object to array	
630	3. Render in next animation frame	
631	– detail.formula: -	
632	<ul> <li>detail.theory: Procedural enemy generation</li> </ul>	
633	- detail.application: Game difficulty balanc-	
634	ing	
635 636	$\texttt{research}_n otes$	
637	• Inspired by grid symmetry in 2023 Shogi tile	
538	patterns	

• Based on known complexity bounds in tiling theory

# References

640

641

642

647

650

651

655

657

658

662

- Tianhao Chen, Weijia Xu, and 1 others. 2024. Pdl: A declarative prompt programming language. *arXiv preprint*.
- Paul F Christiano, Jan Leike, Tom B Brown, Miljan Martic, Shane Legg, and Dario Amodei. 2017. Deep reinforcement learning from human preferences. In Advances in Neural Information Processing Systems, volume 30, pages 4299–4307.
- Zican Dong, Tianyi Tang, Lunyi Li, and Wayne Xin Zhao. 2023. A survey on long text modeling with transformers. *arXiv preprint*.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *arXiv preprint*.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, and et al. 2022. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*.

Sameer Pimparkhede, Mehant Kammakomati, Srikanth G. Tamilselvam, Prince Kumar, Ashok Pon Kumar, and Pushpak Bhattacharyya. 2024. Doccgen: Document-based controlled code generation. *arXiv preprint*. 663

664

665

666

667

668

669

670

671

672

673

674

675

676

- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30, pages 5998–6008.
- Ming Wang, Yuanzhong Liu, and 1 others. 2024. Langgpt: Rethinking structured reusable prompt design framework for llms from the programming language. *arXiv preprint*.