

---

# KNOWLEDGE GRAPH FUSION WITH LARGE LANGUAGE MODELS FOR ACCURATE, EXPLAINABLE MANUFACTURING PROCESS PLANNING

---

**Danny Hoang**

School of Mechanical, Aerospace, and Manufacturing Engineering  
University of Connecticut  
Storrs, Connecticut, USA

**David Gorsich**

The United States Army Combat Capabilities Development Command  
Ground Vehicle Systems Center  
Warren, Michigan, USA

**Matthew P. Castanier**

The United States Army Combat Capabilities Development Command  
Ground Vehicle Systems Center  
Warren, Michigan, USA

**Farhad Imani \***

School of Mechanical, Aerospace, and Manufacturing Engineering  
University of Connecticut  
Storrs, Connecticut, USA

## ABSTRACT

Precision process planning in Computer Numerical Control (CNC) machining demands rapid, context-aware decisions on tool selection, feed-speed pairs, and multi-axis routing, placing immense cognitive and procedural burdens on engineers from design specification through final part inspection. Conventional rule-based computer-aided process planning and knowledge-engineering shells freeze domain know-how into static tables, which become limited when dealing with unseen topologies, novel material states, shifting cost-quality-sustainability weightings, or shop-floor constraints such as tool unavailability and energy caps. Large language models (LLMs) promise flexible, instruction-driven reasoning for tasks such as G-code synthesis to spindle-load queries, but they routinely hallucinate numeric values and provide no provenance. We present Augmented Retrieval Knowledge Network Enhanced Search & Synthesis (ARKNESS), the end-to-end framework that fuses zero-shot Knowledge Graph (KG) construction with retrieval-augmented generation to deliver verifiable, numerically exact answers for CNC process planning. ARKNESS (1) automatically distills heterogeneous machining documents, handbooks, G-code annotations, and vendor datasheets into augmented triple, multi-relational graphs without manual labeling, and (2) couples any on-prem LLM with a retriever that injects the minimal, evidence-linked subgraph needed to answer a query. Benchmarked on 155 industry-curated questions spanning tool sizing, feed-speed optimization, and tolerance diagnostics, a lightweight 3B-parameter Llama-3 augmented by ARKNESS matches GPT-4o accuracy while achieving a +25 Percentage Point (pp) gain in multiple-choice accuracy, +22.4 pp in F1, and 8.1 $\times$  ROUGE-L on open-ended responses. Additionally, by grounding reasoning in precise triples, ARK-

---

\*Corresponding author: farhad.imani@uconn.edu

DISTRIBUTION STATEMENT A. Approved for public release; distribution is unlimited. OPSEC9768

NESS enables smaller models to match or exceed the accuracy of much larger cloud models, reducing numeric hallucinations by 22 pp, and running fully on-premise for privacy-preserving, real-time inference on the shop floor.

**Keywords** Large Language Model · Retrieval Augmented Generation · Knowledge Graph, · Manufacturing Process Planning

## 1 Introduction

Modern production processes, such as precision machining, demands accuracy margins and tolerates virtually no numerical error, imposing substantial mental and operational loads on engineers from initial design specifications to the final part evaluation [1, 2, 3]. In high-stakes sectors, such as aerospace and energy, dimensional tolerances routinely tighten to  $\pm 5 \mu\text{m}$ , while surface-integrity constraints, white-layer thickness, and residual stress, must remain within narrowly defined thresholds [4]. To meet these criteria, machinists must simultaneously optimize a variety of intertwined process parameters, such as tool type, feed rate, spindle speed, and depth of cut, carefully trading off cycle time against surface finish, tool wear, and geometric fidelity [5, 6, 7]. An incorrect decimal drill size or an outdated cutting-speed chart therefore translates directly into scrap, rework or latent defects that may only become apparent after catastrophic failure [8]. Such unplanned downtime erodes an estimated 11% of annual revenue for Fortune Global 500 manufacturers, roughly US \$1.5 trillion lost each year [9].

Computer-Aided-Process-Planning (CAPP) engines and the Computer-Aided Manufacturing (CAM) modules bundled within mainstream Computer-Aided Design (CAD) suites (e.g., Siemens NX, CATIA, and SolidWorks) form automation pipeline that ingests STEP-compliant product models (AP203/242) and emits executable Numerical Control (NC) results [10]. A dedicated CAPP system executes the high-level analysis loop, geometry interrogation, feature recognition, precedence, constraint resolution, raw-stock/fixture/machine assignment, and operation-graph synthesis, before compiling feeds, speeds, and tool-change cycles into ISO6983 or STEP-NC. Embedded CAM modules then refine the middle layers: selecting cutters from vendor libraries, applying parameter presets, generating and simulating toolpath with collision checks and material-removal estimates, and invoking feature-based machining or knowledge-based-engineering rules for common prismatic features. Collectively, these rule-driven frameworks can compress planning lead time by 2–10 $\times$ , yet they share a structural brittleness. Hard-coded heuristics falter when faced with unseen topologies (e.g., self-intersecting free-form pockets), novel material states, shifting cost-quality-sustainability weightings, or shop-floor constraints, such as tool unavailability and energy caps. Moreover, neither tier natively ingests telemetry (e.g., cutting force, vibration, and thermal imagery) generated by machining cells, forcing human experts to intervene, reconcile documentation, and manually retune parameters, an increasingly untenable bottleneck as geometrical complexity and throughput requirements escalate.

Large Language Models (LLMs) have emerged as versatile, instruction-following agents capable of synthesizing domain knowledge, reasoning under uncertainty, and generating executable results [11]. Recent studies demonstrate their ability to synthesize 3-axis drilling paths, to answer contextual queries about spindle loads and feed rates, and to troubleshoot controller alarms with human-level accuracy in descriptive tasks [12]. Šket et al. benchmarked GPT-3.5 and GPT-4 across independent G-code generation, self-interpretation, and error simplification; GPT-4 rendered more correct toolpath yet remained limited to simple drilling operations and required extensive prompt steering [12]. To broaden functionality beyond code synthesis, ChatCNC was introduced, coupling multiple LLM agents with real-time machine telemetry so that an operator can, for example, ask “What was the spindle load at 09:32?”; the system attained 93.3% accuracy on complex production-tracking queries [13]. Nevertheless, both studies expose two systemic bottlenecks: first, dependence on precise user prompts or narrowly structured SQL retrieval, and second, fragile performance when required context is absent from the local database.

Domain-specific fine-tuning offers one mitigation strategy. Rosati et al. raised ROUGE-1 scores of Llama-3 from 0.164 to 0.314 using industrial user manuals [14]. CNC-GPT fine-tuned GPT-3.5 Turbo on machine instructions and trouble-shooting logs, boosting factual correctness from ROUGE-L 0.296 to 0.692 [15]. Yet intensive fine-tuning presents new trade-offs: (1) escalating computational cost as the query-response inventory grows into the thousands; (2) rapid obsolescence that mandates frequent retraining; and (3) data-sovereignty concerns when proprietary part files or production metrics traverse external cloud APIs. Hence, a complementary mechanism is required, one that supplies rich, verifiable context to an LLM without chronically retraining the model or exposing sensitive data.

Knowledge Graphs (KGs), on the other hand, provide a structured, semantics-rich foundation for capturing machining know-how as typed entities (features, tools, operations, machines) and relations (e.g., requiresTool, precedes, causes) [16, 17, 18]. For CAPP, Wang et al. generated a feature-machining KG and proposed a cosine-similarity retrieval that chose industry-valid hole-machining schemes with 0.8450 similarity to expert templates [19]. Guo et al. organized

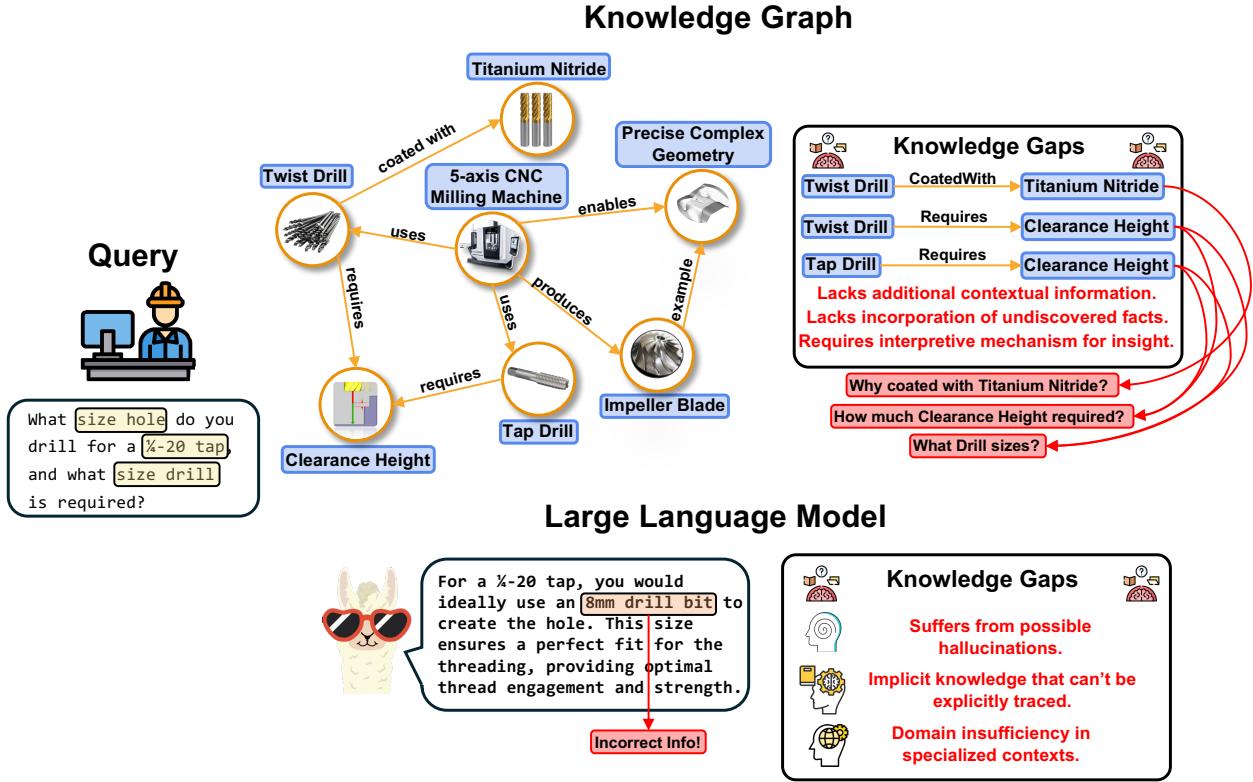


Figure 1: Illustration contrasting a knowledge graph (top) and a large language model (bottom) handling the same tap-drill sizing query. In the knowledge graph section, technical entities (i.e., Twist Drill, Clearance Height, Tap Drill, Titanium Nitride, 5-axis CNC Milling Machine, and Impeller Blade) are connected by explicit relations but can often be missing contextual links and undiscovered facts. The LLM section shows the model’s generated answer (“8 mm drill bit”), annotated to expose imprecision risks, implicit untraceable knowledge, and domain inefficiency in specialized machining.

historical process routes by feature topology and machine capability, trimming part-routing time from 50–80 min to approximately 15 min [20]. Beyond planning, KGs have achieved >90% accuracy in root-cause diagnostics of rotary machinery even under missing sensor data by reasoning over causal chains [21]. Despite these gains, KG deployment still depends on laborious, manual curation: (1) domain experts must annotate source texts, validate triples, and periodically inject new machining knowledge; (2) Simple subject–predicate–object triples also struggle to encode context such as tolerance stack-ups, fixture constraints, or process-chain rationale for nuance components. For example, a given connection between the tool used to create thin walled geometry might omit the dependency between the optimal tool geometry corresponding to the feature size and specific tolerances required for that specific part.

Figure 1 showcases these limitations by contrasting (1) a stand-alone knowledge graph that offers traceable but context-starved reasoning, returning no answer when crucial relations such as material-specific clearance or over-size rules are absent, and (2) a stand-alone LLM that confidently hallucinates an incorrect 8 mm recommendation, revealing the need for a hybrid KG and LLM architecture that is both provenance-aware and context-adaptive. This research introduces Augmented Retrieval Knowledge Network Enhanced Search and Synthesis (ARKNESS), a hybrid framework that combines LLMs with semantically enriched, domain-specific KGs that provides answers that are grounded in validated information ensuring that recommendations for machining process planning queries are both reliable and contextually relevant. In addition, by grounding LLM prompts in these rich subgraphs, ARKNESS can deploy a compact parameter model on-premise with accuracy on par with, or exceeding, much larger cloud LLMs, slashing compute and deployment overhead. The main benefits of this paper are:

1. **A model-agnostic, KG-augmented CAPP assistant.** We present a retrieval-augmented generation pipeline that couples any large language model with a multi-relational machining knowledge graph. A graph-aware prompt compiler injects provenance-linked triples into the context window, yielding numerically precise and explainable process-planning answers.

2. **Self-supervised KG distillation from technical corpora.** A zero-shot, GPT-based entity–relation extractor converts heterogeneous machining documents (e.g., PDF manuals, specification sheets, and NC code comments) into contextualized triples, eliminating the manual curation bottleneck that has constrained previous CAPP-oriented KGs.
3. **Lightweight, on-prem hallucination suppression.** On a 155-query benchmark, ARKNESS paired with a smaller LLM model matches or surpasses much larger cloud LLMs while reducing hallucinations, demonstrating a privacy-preserving path to real-time shop-floor deployment.

The rest of the paper is as follows: Section 2 provides relevant knowledge graph and large language modeling research in manufacturing planning contexts, Section 3 provides the implementation of ARKNESS for automatic graph creation for retrieval and answering, Section 4 describes the experimental setup, Section 5 presents the experimental results, and Section 6 concludes the paper.

## 2 Research Background

### 2.1 Large Language Models in Machining

LLMs are rapidly transitioning from research curiosities to core enablers of AI-driven production, thanks to their capacity to parse free-form instructions, fuse heterogeneous context, and return actionable, domain-specific guidance. A recent survey by Li et al. categorized their early penetration into the manufacturing stack—spanning generative CAD modeling, bio-process recipe design, robot path planning, and vision-based quality control [22]. Within subtractive manufacturing, the most direct application is natural-language-to-G-code translation. Šket et al. evaluated commercial ChatGPT models for 3-axis G-code generation testing GPT-3.5 and GPT-4 in three phases including independent G-code generation, interpretation of the generated G-code, and detecting and simplifying errors [12]. Their results showed promise of implementing LLMs for G-code generation with GPT-4 producing more correct toolpaths but is severely limited to simple operations such as drilling. Additionally, their method relied heavily on user input to align the LLMs for generation which can lead to increased downtime if deployed in manufacturing environments. Jeon et al. expanded upon just G-code generation by developing ChatCNC that integrated various LLM agents with real-time CNC machining data [13]. This allowed users receive context-aware answers regarding the status of their 3-axis CNC machines such as spindle load at specific instances or times; their method achieved an accuracy of 93.3% in queries requiring complex data inference such as production tracking showcasing applications in analyzing data recorded in the manufacturing pipeline. Despite the high accuracy, the authors acknowledged that the model resulted in failures when encountered with missing context or the inability to retrieve information from their database. This suggests a more in depth search beyond conventional retrieval methods such as defined SQL database structures is required to supplement existing knowledge gaps.

Beyond prompting and traditional retrieving techniques, researchers have also introduced domain-specific fine-tuning to better align a LLMs responses with domain-specific knowledge. Rosati et al. finetuned Llama 3 for industrial applications achieving uplifts in average ROUGE-1 F1-score from 0.164 to 0.314 when trained on user manuals of a 360-degree camera [14]. Wang et al. finetuned GPT-3.5 for air craft maintenance outperforming general GPT-3.5 and its upgraded counterpart GPT-4.0 [23]. In areas specific to machining, Soundararajan et al. developed CNC-GPT aimed for on-site CNC operator assistance by finetuning GPT-3.5 Turbo on machine-specific data, operational instructions, and troubleshooting assertions [15]. Results showed an increase in factual correctness when using ROGUE-L scores from 0.296 to 0.692 before and after finetuning, respectively. Despite the validity of their method, there are some challenges that might prevent further deployment on factory floors. One such issue is the amount of resources and training time required as the authors only trained around thirty query-response examples. In real world scenarios, there might be hundreds to thousands of scenarios and responses which can severely limit scalability during training. Moreover, as new scenarios continually emerge, the system must be frequently retrained to stay current which can further strain computational resources and complicates maintenance. There is also the challenge of security regarding production processes and information contained within. The authors used commercial LLMs which operate on external, cloud-based platforms where data flows and storage may raise concerns about data privacy and intellectual property protection. This reliance on third-party systems increases the risk of unauthorized access or data leakage, making it imperative to implement solutions that explore on-premise solutions to ensure that proprietary manufacturing data remains strictly confidential.

### 2.2 Knowledge Graphs in Machining

The ability to reuse and implement machining knowledge in a structured way provides not only enhanced consistency across planning and decision making but also ensures existing knowledge gaps are mitigated. Xiao et al. reviewed

how computer-aided process planning involving knowledge graphs can benefit from reduced labor costs, shortened production cycles, and more intelligent use of existing information [16]. The authors analyzed key steps of implementing knowledge graphs from process knowledge representation to process knowledge graph construction and validation, showcasing how these methods can overcome traditional CAPP by reducing excessive manual interventions, and increasing flexibility and generalization. This implementation of process knowledge was shown in Wang et al. where they constructed a process knowledge graph for feature-based machining to automate machining scheme selection [19]. Through use of an improved cosine-similarity formula for machining scheme selection, they achieved a similarity score of 0.8450 closely matching existing mature schemes implemented in industry for a typical shell part composed of 6 holes. By using their method, Wang et al. argued that the recommended machining steps would reduce tool load and increase both part quality and machining safety. In a similar study, Guo et al. created a knowledge graph for process route reuse by organizing historical process plans using part feature topology and machine capabilities [20]. When introducing a new part, the authors determine the process route that best aligns with its feature topology through a similarity check with existing process routes. Their case study on a shaft part resulted in completing the overall machining process route in about 15 minutes, down from the typical 50 to 80 minutes, significantly enhancing efficiency.

Another major application of machining knowledge graphs involves fault diagnosis and maintenance of CNC equipment. Manufacturing systems generate heterogeneous data from sources such as sensors, logs, and maintenance reports that if left in isolation or silos can be difficult for diagnosis if problems arise. Knowledge graphs offer a unified representation by connecting physical components, signals, and failure modes in a network of cause-effect and part-whole relationships. Qiu et al. tackled this "data island" problem by constructing a multi-layered knowledge graph that integrated sensor data and domain knowledge, enabling the automatic identification of health changes in the X-axis ball screw drive system on a vertical milling tool (model XHK-5140) over 201 days, while similarity-based reasoning over the graph quantified the deviation from a healthy baseline [24]. Building on this idea, Cai et al. introduced a multi-level fault diagnosis KG for rotating machinery combining hierarchical knowledge of subsystem states with Bayesian reasoning, offering probabilistic inference across the graph to pinpoint root causes [21]. Notably, even with missing sensor outputs, their method achieved 91.1% diagnostic accuracy by leveraging the relationships between related symptoms and outperformed traditional rule-based fault diagnosis.

Despite the considerable progress in implementing knowledge graphs, there remain challenges in deploying them at scale for machining and industry. These aforementioned methods continue to depend heavily on extensive manual curation. The process of extracting, validating, and continuously updating the domain knowledge requires significant human knowledge and intervention, making it difficult to scale rapidly or adapt quickly to new information. This labor-intensive approach not only increases the risk of inconsistencies or omissions but also impedes real-time responsiveness in dynamic manufacturing environments. Furthermore, the conventional reliance on simple triple-based representations constrains the ability to capture the full contextual richness necessary for nuanced decision-making in machining processes. Such representations often fall short in encoding the detailed circumstances, rationale, and intricate relationships that underlie complex manufacturing operations. As a result, they may fail to fully express how various machining operations interrelate or why specific process choices are chosen. This would thus lead to a superficial understanding of the underlying mechanics and process dependencies. These limitations highlight the need for more flexible and richly contextualized approaches in knowledge representation to support the evolving demands of modern machining applications.

### 3 Research Methodology

This section details the two main components of the ARKNESS framework with 1) Knowledge graph construction and 2) Graph transversal, retrieval and LLM response. Figure 2 provides the general overview of the two components.

#### 3.1 Knowledge Graph

A knowledge graph can be considered as a set of tuples  $\mathcal{G} = \{(v, r, w)\}$  where  $v$  and  $w$  are vertices (entities) from a set of vertices  $\mathcal{V}$  and  $r$  is an edge (relationship) from a set of relationships  $\mathcal{R}$ . The set of tuples is commonly known as triples denoted as  $(\text{head}, \text{relation}, \text{tail})$  representing  $v$ ,  $r$ , and  $w$ , respectively.

#### 3.2 Automated Graph Construction

The knowledge graph construction for ARKNESS is mainly done using textual data from documents relevant to the users chosen domain. Due to the heterogeneous nature of document types and file formats (e.g., .docx, .pdf, .pptx), preprocessing was first done to extract the textual information from these different sources. This was achieved through

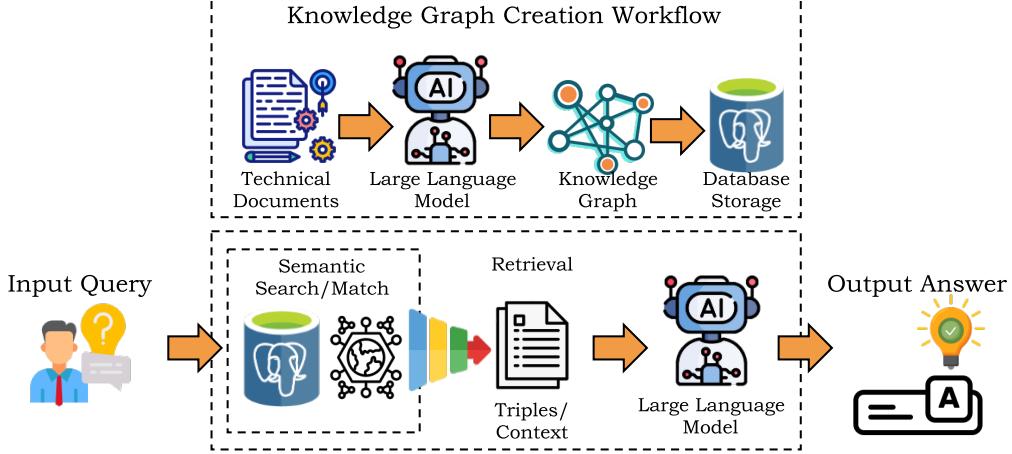


Figure 2: Overview of the framework with incorporating knowledge graph with large language model, illustrating the process from graph creation and storage to semantic search, match, and retrieval, and finally to language model response.

the Docling [25] Python package which efficiently parses through the documents and extracts the raw text subsequently exporting the content as Markdown files for further processing.

Document splitting was done through parsing out each individual paragraph. Given a document with  $p$  individual paragraphs, each paragraph was given to a LLM to extract entities and relations forming  $T$  triples. Specific instructions were given to guide the LLM to output structured information given below:

```
-Goal-
Extract structured triples directly from the input text in the following format:

ENTITY_1, RELATIONSHIP_TYPE, ENTITY_2, "RELATIONSHIP_DESCRIPTION"

-Steps-
1. Read the input text carefully to identify:
   - Key entities: Concepts, systems, technologies, or processes central to the text.
   - Relationships: Clear actions or connections described in the text that link two entities.
   - Descriptions: Verbatim or paraphrased descriptions from the text that explain the relationship.

2. For each relationship, construct a triple:
   - ENTITY_1: The primary concept or entity initiating the relationship.
   - RELATIONSHIP_TYPE: The action or connection type as described in the text.
   - ENTITY_2: The target concept or entity affected by ENTITY_1.
   - RELATIONSHIP_DESCRIPTION: A concise description of the relationship directly sourced from the input text.

3. Each triple must be clear and in this format:
   ENTITY_1, RELATIONSHIP_TYPE, ENTITY_2, "RELATIONSHIP_DESCRIPTION"

4. Use the original text verbatim where possible for the description, ensuring accuracy.
   Avoid adding external interpretations or explanations.
```

This process was repeated for each document chosen creating  $G_D$  subgraphs. The purpose of including the 'Relationship Description' in the triplet information is to retain the relevant contextual information by which the triplet was created. This removes ambiguity and often missing information found in traditional knowledge graphs. All subgraphs were then stored into respective text files which were then combined together for further processing.

A hierarchical database was then created using PostgreSQL to efficiently store and retrieve entities, relations, and their associated context. Specifically, Python and the Psycopg library was used to interface with the database where during initialization three primary tables are created: (1) subjects table: stores unique entity names with a serial primary key, (2) relations table: associates each subject with its corresponding relationships where each entry references a subject's unique identifier, ensuring relationships are correctly linked, and (3) objects table: stores ending triplet entities along with corresponding contextual information. This structure can be also represented as follows. Let  $S$  be the set of subjects (entities). For each subject  $s \in S$ , let  $R(s)$  be the set of relations associated with  $s$ . For each relation  $r \in R(s)$ , let  $O(r)$  be the set of ending objects linked to  $r$ , where each object is represented as a tuple  $(v, c)$  with  $v$  representing connected entities, and  $c$  representing the corresponding context. The entire data base is thus defined by the set:

$$\mathcal{D} = \{\langle s, r, (v, c) \rangle | s \in S, r \in R, (v, c) \in O(r)\} \quad (1)$$

with a single triple entry represented as:

$$\langle s, r, (v, c) \rangle \quad (2)$$

### 3.3 Graph Retrieval and Transversal

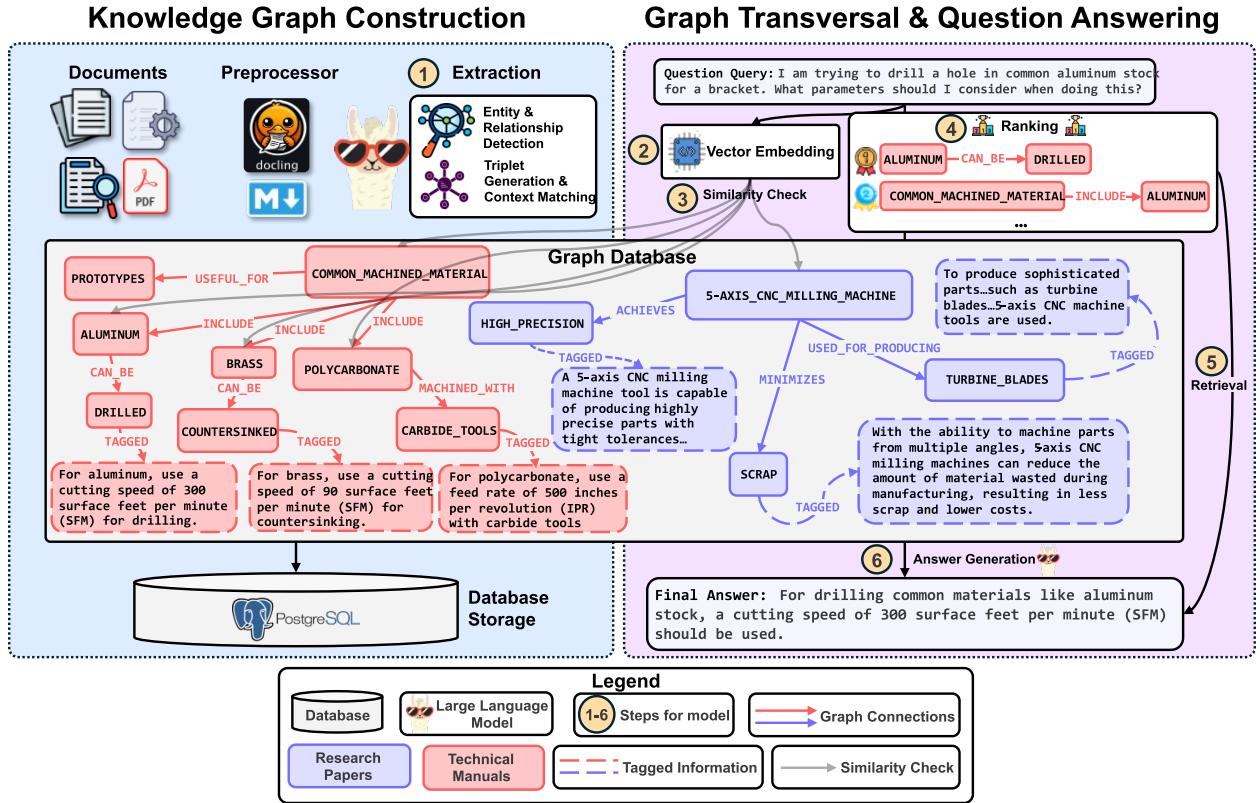


Figure 3: Overview of the knowledge graph construction and graph transversal for user queries. The left side depicts the automated graph construction process from document selection, preprocessing, entity-relationship extraction, triplet generation and context matching, and database storage. The right side depicts knowledge graph retrieval given a user query by embedding the query to get a vector representation, performing a similarity check over the triples in the graph, and ranking the most relevant triples to give to the large language model for answer generation.

Figure 3 provides an indepth illustration of how triplets are compared against a query and retrieved to be given to a LLM. Once the database is created and stored, the information contained within is incorporated to provide additional context and information to anchor subsequent questions. This is achieved through retrieving the relevant triplets and context from the knowledge graph. The retrieval process begins by first encoding a given natural language query, or question, into a vector  $q \in \mathbb{R}^d$ . This is achieved using a chosen semantic embedding model  $f : \text{Text} \rightarrow \mathbb{R}^d$ . Similarly, each triplet  $t_i$  from the database is encoded into embedding  $e_i$  using the same semantic embedding model where  $f : t_i \rightarrow \mathbb{R}^d$ .

The similarity between the query and each triple is then calculated using cosine-similarity according to:

$$\text{sim}(q, e_i) = \frac{q \cdot e_i}{\|q\| \|e_i\|} \quad (3)$$

This similarity score provides an evaluate of how semantically similar the query given is to each of triplet in the knowledge graph allow one to ascertain its relevance.

Given the potentially massive scale of knowledge graphs, which can reach hundred of thousands or more connections, a candidate pool  $C$  is constructed by selecting indices corresponding to a chosen a top- $K$  triples with the highest similarity scores  $s_i = \text{sim}(q, e_i)$ . Having only the highest similarities chosen allows the most semantically relevant triples to be selected for further processing. The candidate pool is thus defined as:

$$C = \{i \in \{1, 2, \dots, K\}\} \text{ with } s_{i_1} \geq s_{i_2} \geq \dots \geq s_{i_K} \quad (4)$$

As mentioned previously, because a knowledge graph can contain hundred of thousands or even millions of connections, relying solely on the top- $K$  triples could overlook additional relevant triples connected to them. To explore the extended neighborhood of these triples, we transverse through the graph by performing a beam search expansion. Let  $E(i)$  denote the set of adjacent embedded triples indexed by  $i$  (i.e., those sharing a node with triple  $i$ ). We conduct a beam search recursively over a chosen maximum expansion depth  $d_{max}$ . Here higher depths allows one to explore deeper through the graph allowing more information to be retrieved. For an initial depth  $d = 0$ , we first set the initial candidates as  $C_0 = C$ . For each subsequent depth  $d \geq 1$  we find additional candidates using:

$$C_d = \bigcup_{i \in C_{d-1}} \text{Top-}b\left(\{j \in E(i) \mid \text{sim}(q, e_j)\}\right) \quad (5)$$

where  $\text{Top-}b(\cdot)$  selects the chosen number of  $b$  candidates, also known as beam width, with the higher similarity scores relative to the embedded query  $q$ . Here, beam width is a chosen parameter where higher values of  $b$  allows the model to explore a wider set of candidates to expand from for each step. This can potentially capture more semantically relevant triples missed in the initial top- $K$ . Conversely, lower values of  $b$  restricts the search to the more promising candidates which reduces computational overhead at the expense of possibly missing relevant connections that lie outside the top  $b$  candidates. In order to avoid re-processing triples, a set of visited triples  $E_v$  is maintained such that the recursive expansion only considers new nodes defined as:

$$C_{d+1} = \bigcup_{i \in C_d \setminus V} \text{Top-}b\left(\{j \in E(i)\}\right) \quad (6)$$

After recursively expanding through the chosen depth and beam size, the final set of triples retrieved is defined by:

$$C^* = \bigcup_{d=0}^{d_{max}} C_d \quad (7)$$

As a final step after the final set of candidate triples are chosen, the corresponding context  $c$  from the database is retrieved as well. This creates a set  $\mathcal{C}$  of relevant information:

$$\mathcal{C} = \{\langle s_i, r_i, (o_i, c_i) \rangle : i \in I\} \quad (8)$$

### 3.4 Large Language Model Generation

The last step of answering the query  $q$  involves prompt construction where the information retrieved from the knowledge graph is given to the large language model to help its answering. This is done by defining a set of system instructions that designate that the LLM is design to answer questions incorporate both its own knowledge and the given retrieved knowledge. For example, we define the system instructions as:

'You are designed to help answer questions using retrieved knowledge. Not all knowledge given  $I_{sys}$  = need to be used but focus on the most important information. Remember this knowledge, if there is any, to help your decision making.' (9)

Then the final prompt  $P$  given to the LLM is constructed by concatenating the system instructions, query, and retrieved knowledge graph information:

$$P(q, \mathcal{C}) = I_{sys} \oplus \mathcal{C} \oplus q \quad (10)$$

The final answer  $A_{LLM}$  is then generated by the LLM:

$$A_{LLM} = \text{LLM}(P(q, \mathcal{C})) \quad (11)$$

## 4 Experimental Design

### 4.1 Knowledge Graph

To test the efficacy of ARKNESS, we first created our machining knowledge graph composed of research and technical documents spanning 5-axis milling capabilities, industrial robotic integration in CNC machining, sustainability in CNC machining, CNC programming and process planning, and fundamentals of CNC machining. By exposing ARKNESS to a spectrum of complexity levels, domain-specific terminologies, and the assembled documents enables a comprehensive evaluation of the model's CNC machining knowledge. Each document was processed using the automated graph-construction pipeline detailed in Section 3.2 using GPT-4o [26]. An example regarding the output given the text information is given in the following box [27]. Descriptions of each documents chosen are listed in Table 1 along with the total number of entities, relations, and triples after processing. In total, the full knowledge graph contains 4329 triples, 6659 unique entities, and 1251 unique relations providing a relatively large graph to test to test ARKNESS' capacity robust inference, and accurate knowledge retrieval across diverse CNC-machining scenarios. An example of the graph-construction output given text input is given below:

<b>Input:</b>
A 5-axis CNC milling machine tool is a sophisticated piece of equipment used for precision machining of complex parts with high accuracy and consistency. It has five axes of motion that can move the cutting tool in various directions to create intricate shapes and contours. The machine is designed to perform a wide range of operations, including milling, drilling, cutting, and engraving, using multiple axes of movement. A 5-axis CNC milling machine tool can create highly complex geometries using five axes of motion during machining operations that which can be difficult or impossible to be implemented by using traditional 3-axis machine.
<b>Output:</b>
5-AXIS CNC MILLING MACHINE TOOL, is used for, PRECISION MACHINING, "A 5-axis CNC milling machine tool is a sophisticated piece of equipment used for precision machining of complex parts with high accuracy and consistency."
5-AXIS CNC MILLING MACHINE TOOL, has, FIVE AXES OF MOTION, "It has five axes of motion that can move the cutting tool in various directions to create intricate shapes and contours."
5-AXIS CNC MILLING MACHINE TOOL, performs, "MILLING, DRILLING, CUTTING, AND ENGRAVING", "The machine is designed to perform a wide range of operations, including milling, drilling, cutting, and engraving, using multiple axes of movement."
5-AXIS CNC MILLING MACHINE TOOL, creates, HIGHLY COMPLEX GEOMETRIES, "A 5-axis CNC milling machine tool can create highly complex geometries using five axes of motion during machining operations that which can be difficult or impossible to be implemented by using traditional 3-axis machine."
TRADITIONAL 3-AXIS MACHINE, limits, GEOMETRY COMPLEXITY, "Machining operations that can be difficult or impossible to be implemented by using traditional 3-axis machine."

### 4.2 Models

To rigorously assess how model choice shapes our framework's performance, we evaluate a diverse set of both open-source and closed-source large language models across multiple parameter scales. For open sources models, Llama 3.2 3B Instruct [36], LLama 3.1 8B Instruct [37], and Qwen 2.5 7B Instruct [38] were chosen due to their ease of access and relatively low computational requirements; here B represents the number of parameters in billions. By spanning 3B, 7B, and 8B parameter tiers, we can isolate how model capacity and design choices interact with our knowledge graph augmentation. Evaluating these models across a broad spectrum of parameter scales enables us to quantify how the integration of supplementary knowledge graph information influences overall performance. To establish an upper bound on attainable performance, we also benchmark against state-of-the-art closed sourced models including GPT-4o [26], its smaller GPT-4o-mini, Gemini 2.0 Flash [39], and Gemini 2.0 Flash-Lite [40]. These comprehensive models enables us to quantify precisely how supplemental structured knowledge closes the gap between open source baselines and leading proprietary offerings.

Table 1: Overview of source documents and their corresponding knowledge graph sizes.

Document	Description	Knowledge Graph Size
A Review in Capabilities and Challenges of 5-Axis CNC Milling Machine Tool Operations [27]	This review surveys recent advances and challenges in 5-axis CNC milling, including error modeling and compensation, toolpath and process optimization, virtual machining systems, tool wear and temperature prediction, and sustainability considerations.	# of triples: 464 # of unique entities: 654 # of unique relations: 199
Design and development of a CNC machining process knowledge base using cloud technology [28]	This paper presents a cloud-based CNC machining process knowledge base that maps STEP-NC to an OWL ontology and leverages Hadoop’s HBase for scalable storage, MapReduce driven querying, and SWRL-based reasoning to enable intelligent, high-throughput process planning.	# of triples: 559 # of unique entities: 828 # of unique relations: 291
Exploring the Application of Industrial Robots in CNC Machining [29]	This paper explores the deployment of industrial robots in CNC machining. It details robot selection, workflow and control program design, joint and transmission mechanisms, and drive system choices to automate loading/unloading and enhance operational flexibility, precision, and efficiency.	# of triples: 98 # of unique entities: 151 # of unique relations: 57
Fundamentals of CNC Machining [30]	This guide provides a practical introduction to CNC machining fundamentals covering shop safety, tooling, coordinate systems, programming and operation of mills and lathes, 2D/3D toolpaths, workholding examples, and best practices for prototype and short-run production.	# of triples: 1902 # of unique entities: 2966 # of unique relations: 512
Innovative Approaches to Sustainable CNC Machining: A Machine Learning Perspective on Energy Optimization [31]	This paper develops a machine-learning framework for sustainable 5-axis CNC milling by combining per-axis power monitoring with Taguchi experimental design and tree-based regression to model and predict energy consumption.	# of triples: 278 # of unique entities: 452 # of unique relations: 183
Integration of Taguchi and PROMETHEE for CNC Milling Machining Parameter Optimization [32]	This paper integrates Taguchi’s orthogonal-array experimental design with the PROMETHEE multi-criteria ranking method to optimize spindle speed, feed rate, and depth of cut for CNC milling of AA6061.	# of triples: 227 # of unique entities: 366 # of unique relations: 74
Research on CNC programming and machining process based on CAD/CAM technology [33]	This paper introduces a CAD/CAM-based CNC programming framework that uses Bayesian process-skeleton mapping to link part features with reusable macro-processes and applies residual-height trajectory generation and nonlinear-error compensation to optimize multiaxis toolpaths.	# of triples: 269 # of unique entities: 443 # of unique relations: 167
Review on Design Research in CNC Machine Tools Based on Energy Consumption [34]	This review surveys global advances in modeling, designing, and evaluating CNC machine tools from an energy-consumption perspective.	# of triples: 270 # of unique entities: 470 # of unique relations: 140
Robotical Automation in CNC Machine Tools, A Review [35]	This review surveys robotics-driven advances in CNC machining. Automated material handling and tool changing to adaptive control, quality inspection, data analytics, and collaborative robots are evaluated for their impacts on efficiency, precision, and safety.	# of triples: 262 # of unique entities: 380 # of unique relations: 129

### 4.3 Question Category

Following knowledge graph construction, we devised two question formats namely, multiple choice and open ended to evaluate each model’s capabilities. Multiple choice questions present a controlled benchmark for retrieval precision, enabling objective measurement and direct comparison across models, while open-ended questions simulate real-world case studies by challenging models to perform generative synthesis and coherently integrate the retrieved information into comprehensive answers. For each question format, two categories of questions were created namely: content

specific and machining specific. Content specific questions refer to the information written in the documents that do not require quantitative analysis. Machining specific questions refer to questions that require quantitative precision and decision-making based on numerical parameters. An example of each question in multiple choice format is given below:

### Content Specific Question

Which statement best describes a key distinction between 3-axis and 5-axis CNC milling machines?

- A. 3-axis machines are used for metals, while 5-axis machines are used exclusively for plastics.
- B. 5-axis machines include tilt and rotation of the workpiece or tool, in addition to X, Y, and Z motion.
- C. 3-axis machines are larger and require more floorspace than 5-axis machines.
- D. 5-axis machines do not allow any vertical movement while 3-axis machines do.[4pt]

### Machining Specific Question

For drilling operations on stainless steel (303), what feed (in inches per revolution) is recommended?

- A. 0.0005 in/rev
- B. 0.0010 in/rev
- C. 0.0020 in/rev
- D. 0.0030 in/rev

GPT-4o was implemented to generate each question from the documents chosen with a total of 65 content specific multiple choice, 45 machining specific multiple choice, 104 content specific open ended, and 45 machining specific open ended.

## 5 Experimental Results

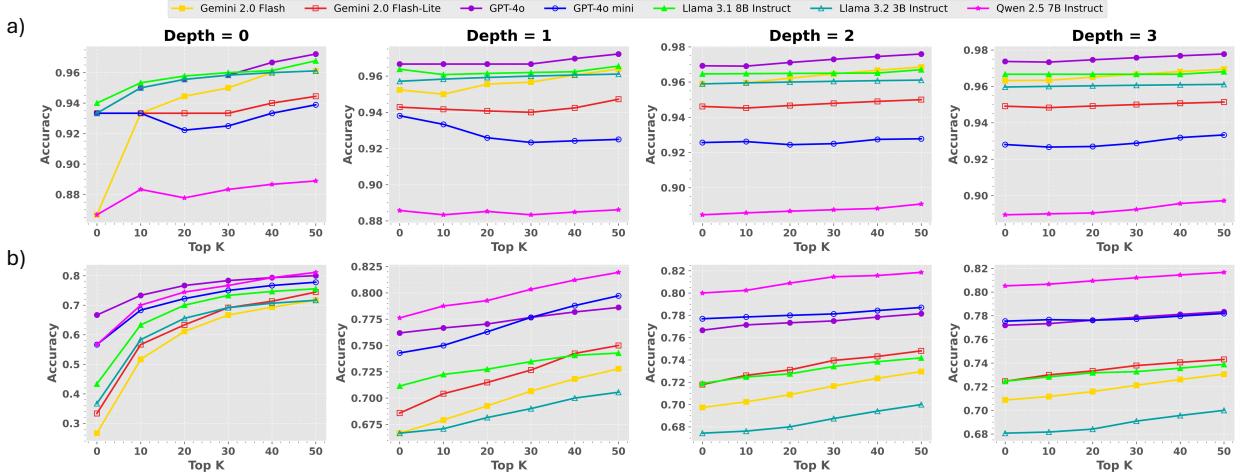


Figure 4: Graph of large language model model accuracy as a function of graph traversal depth and Top K triple retrieval for a) content specific queries and b) machining specific queries.

ARKNESS is implemented using Python, PyTorch, and respective APIs for each of the chosen models. We initially explore how knowledge graph augmentation can enhance the performance of large language models using controlled graph traversals as defined by Eq. 5. Figure 4 shows experiments that vary the depth of these traversals and measure the effect on the accuracy of the LLM answers. In these experiments, we also focus on selecting the Top K highest scoring triples along with their corresponding context. Figure 4a presents results with multiple choice questions that were designated content specific. As shown, as we increase the number of Top K triples retrieved from the graph to give to the LLMs the general trend is that performance increases. This trend holds true across most depths tested as there is an exception where GPT-4o mini at depth one experiences a decrease in performance, possibly because the additional graph-derived information is redundant or less informative, overwhelming the model's initial context. There is a notable

increase for Gemini 2.0 Flash at depth 0 between no information given (Top K of 0) to Top K of 10 increasing from an accuracy of 0.867 to 0.933. It can also be observed that as the depth increases, the performance across the retrieved Top K stabilizes, indicating that for these types of questions a depth of 0 is most optimal since additional information does not yield further benefits. This stabilization implies that the initial retrieved context is already sufficiently rich, and incorporating extra graph-derived details beyond depth zero fails to enhance model accuracy and may even introduce superfluous noise. Compared to content-specific questions, we observe a significantly greater trend for quantitative machining specific questions, with performance improvements manifesting across all depths and Top K values as shown in Figure 4b). This indicates that additional graph-derived context plays a crucial role in accurately identifying the correct quantitative values for machine-specific queries, with stabilization around a depth of 2 demonstrating that these more challenging questions benefit from a deeper graph traversal to capture the necessary relevant information.

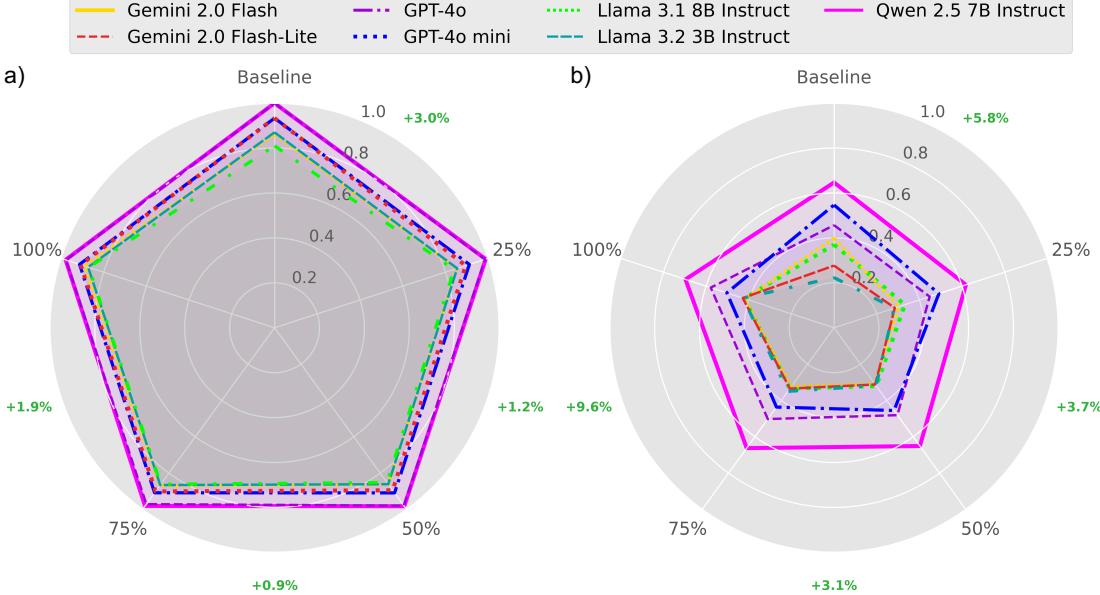


Figure 5: Radar plots of average large language model accuracy performance across knowledge graph completeness levels ranging from baseline (0%) to full graph (100%) for a) content Specific and b) machining specific multiple choice queries.

Figure 5 showcases our model under conditions of incomplete knowledge graphs to emulate real-world scenarios, where the available data may be partial or sparsely connected. This approach not only tests the robustness of the graph augmentation mechanism but also helps us understand how missing nodes and relations affect the retrieval of relevant information and, ultimately, the accuracy of the LLM responses. We again tested the same aforementioned LLMs under varying graph sizes, specifically, baseline with no graph, 25%, 50%, 75%, and full graph of 100%. We randomly dropped triples to achieve the chosen graph sizes and each model was tested 10 times for each respective graph size with the average reported on the plot and the green percent increase indicating the maximum increase across all models tested when increasing the graph size. Figure 5a presents the performance metrics for content-specific questions across various LLMs. The results indicate only a modest improvement when additional graph context is provided, suggesting that these questions can largely be addressed with the baseline contextual information available to each model. Notably, the maximum performance gain, an absolute increase of approximately 3%, is observed when the graph size is augmented from the baseline to 25%, demonstrating that while supplementary retrieved data offers some benefit, its impact remains relatively limited for content-specific queries. Figure 5b illustrates how expanding the graph size affects performance on the more specialized quantitative machining specific questions. Unlike the content specific queries in Figure 5a, these questions exhibit more pronounced gains as the graph grows beyond the baseline level indicated by the greater maximum increase across the models. When comparing closed source models with open sourced, GPT-4o effectively leverages the additional graph context to improve performance by 12.1% at a graph size of 100% compared to baseline. Additionally, Gemini 2.0 Flash-Lite achieves a relatively high increase of 15.1 % compared to baseline at a graph size of 100%. In contrast, the open source Llama 3.2 3B Instruct shows an even greater relative gain of 20.3% at graph size of 100% compared to baseline, suggesting that these models, which may have a lower baseline performance, can be significantly boosted by enhanced external context. Notably, the significant improvements seen in the smaller open source models imply that in manufacturing environments, such as the factory floor, these

models may be particularly effective when enhanced with additional contextual information. Their lower computational overhead and faster inference speeds make them especially suited for deployment in real-time industrial settings where cost efficiency and rapid responsiveness are critical.

Table 2: Benchmark accuracy and F1-score performance between baseline and with knowledge graph for each large language model. Absolute increases (or no change) are in **green**, and decreases are in **red**.

Model	File	Baseline		Knowledge Graph	
		Accuracy	F1 Score	Accuracy	F1-Score
Gemini 2.0 Flash	Content Specific	0.867	0.535	0.933 ( <b>+0.066</b> )	0.569 ( <b>+0.034</b> )
	Machining Specific	0.267	0.152	0.517 ( <b>+0.250</b> )	0.224 ( <b>+0.072</b> )
Gemini 2.0 Flash-Lite	Content Specific	0.933	0.888	0.933 ( <b>0.000</b> )	0.688 ( <b>-0.200</b> )
	Machining Specific	0.333	0.277	0.567 ( <b>+0.234</b> )	0.338 ( <b>+0.061</b> )
GPT-4o	Content Specific	0.933	0.716	0.950 ( <b>+0.017</b> )	0.719 ( <b>+0.003</b> )
	Machining Specific	0.667	0.464	0.733 ( <b>+0.066</b> )	0.573 ( <b>+0.109</b> )
GPT-4o mini	Content Specific	0.933	0.906	0.933 ( <b>0.000</b> )	0.901 ( <b>-0.005</b> )
	Machining Specific	0.567	0.391	0.683 ( <b>+0.116</b> )	0.549 ( <b>+0.158</b> )
Llama 3.1 8B Instruct	Content Specific	0.940	0.925	0.953 ( <b>+0.013</b> )	0.940 ( <b>+0.015</b> )
	Machining Specific	0.433	0.244	0.633 ( <b>+0.200</b> )	0.404 ( <b>+0.160</b> )
Llama 3.2 3B Instruct	Content Specific	0.933	0.718	0.950 ( <b>+0.017</b> )	0.712 ( <b>-0.006</b> )
	Machining Specific	0.367	0.319	0.583 ( <b>+0.216</b> )	0.543 ( <b>+0.224</b> )
Qwen 2.5 7B Instruct	Content Specific	0.867	0.580	0.883 ( <b>+0.016</b> )	0.614 ( <b>+0.034</b> )
	Machining Specific	0.567	0.371	0.700 ( <b>+0.133</b> )	0.463 ( <b>+0.092</b> )

Table 2 provides a benchmark comparison between baseline model performance and results with a minimal viable integration of knowledge graph information (i.e., a chosen Top K of 10 and depth of 0) for both multiple choice datasets. Results for both accuracy and F1-score are reported with the absolute increases, or no increase, between baseline and knowledge graph augmented outputs highlighted in green and any declines shown in red. As shown in the table, for content specific questions where baseline performance was already high, the improvements are very modest and range from 0.013 to 0.066 for accuracy. However, even this relatively small uplift indicates that the inclusion of the knowledge graph imparts valuable domain-specific insights that further refine the responses. On the other hand, for the quantitative machining specific questions the results are far more substantial. For example, Gemini 2.0 Flash-Lite achieves an uplift of 0.250 in accuracy and 0.072 in F1-score. Similarly, magnitudes of improvement is shown with Llama 3.2 3B Instruct and Llama 3.1 8B Instruct having an uplift of 0.216 and 0.200 in accuracy, respectively. These improvements provide evidence that the structured, domain-specific context provided by the knowledge graph effectively addresses the deficiencies each LLM’s internal representations. It should be noted that there is a notable decrease in F1-score of 0.200 for Gemini 2.0 Flash-Lite when testing on content specific questions. This decline may be attributed to the additional context provided by the knowledge graph, which, while generally beneficial, can increase the complexity of the input and lead to confusion in the model’s internal representations and final decision making. In summary, despite some localized declines, likely due to increased input complexity, the integration of knowledge graph information consistently enhances performance across most models, providing significant improvements in domain-specific accuracy and reliability across both content specific and quantitative machining specific queries, and thereby demonstrating its critical value for augmenting LLM outputs in technical applications.

Our next experiment examines the ability of LLMs to respond to open-ended questions versus multiple-choice questions. In this section, we seek to understand how the models perform when required to generate answers without predefined options, compared to when they can select from a limited set of potential responses. This approach allows us to assess the effect of external knowledge augmentation in more unconstrained scenarios. Figure 6 provides the average results across 10 runs for each of the datasets for four different metrics: semantic similarity, ROUGE-1, ROUGE-2, and ROUGE-L. Here, semantic similarity evaluates the closeness in meaning between the generated answer and the reference response by comparing their vector embeddings, thereby capturing the overall semantic content beyond mere word matching. ROUGE-1 measures the overlap of individual words (unigrams) between the generated text and the grounded answer, serving as a baseline indicator of lexical similarity. ROUGE-2 extends this comparison to pairs of consecutive words (bigrams), offering insights into the consistency of short phrases. Finally, ROUGE-L assesses the longest common subsequence between the generated and reference texts, highlighting the structural similarity and the preservation of narrative flow. As shown in Figure 6a), for the content specific questions with the additional information from the knowledge graph, the overall scores are higher for most of the models across all metrics. This is especially true

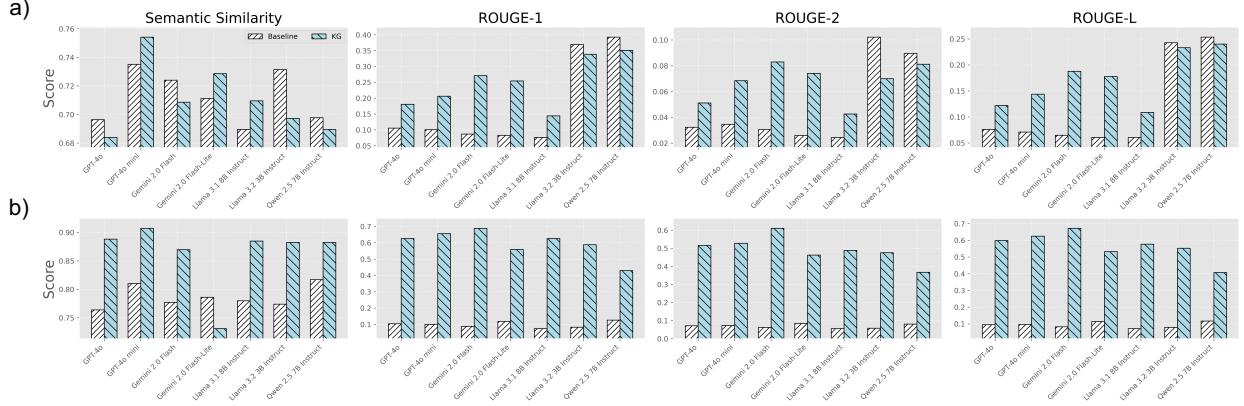


Figure 6: Average performance over 10 runs for evaluated large language models, comparing baseline outputs with those augmented by knowledge graph data. a) Shows performance for content specific questions, and b) depicts results for machining specific questions across semantic similarity, ROUGE-1, ROUGE-2, and ROUGE-L metrics.

Table 3: Open ended responses from open source models Llama 3.2 3B Instruct, Llama 3.1 8B Instruct, and Qwen 2.5 7B Instruct regarding drill sizes for CNC machining. Grounded quantitative answers are highlighted in yellow, baseline values in red, and knowledge graph augmented values in green.

Model / Question	What is the decimal equivalent of Drill Size 82?	What is the decimal equivalent of Drill Size 84?	What is the decimal equivalent of Drill Size 89?
<b>Grounded Answer</b>	The decimal equivalent of Drill Size 82 is <b>0.0125 inches</b> .	The decimal equivalent of Drill Size 84 is <b>0.0115 inches</b> .	The decimal equivalent of Drill Size 89 is <b>0.0091 inches</b> .
<b>Llama 3.1 8B Instruct</b>	<b>Baseline:</b> We find that Drill Size 82 is actually equivalent to <b>0.082 inches</b> . This is because the number 82 represents the 82nd drill size in the series, and the decimal equivalent is calculated using a formula that takes into account the drill size and the series.	<b>Baseline:</b> The decimal equivalent of drill size 84 is <b>0.84 inches</b> , which can also be represented as a fraction, 84/64.	<b>Baseline:</b> The decimal equivalent of Drill Size 89 is approximately <b>0.139 inches</b> .
	<b>Knowledge Graph:</b> The decimal equivalent of Drill Size 82 is <b>.0125 inches</b> .	<b>Knowledge Graph:</b> The decimal equivalent of Drill Size 84 is <b>.0115 inches</b> .	<b>Knowledge Graph:</b> The decimal equivalent of Drill Size 89 is <b>.0091 inches</b> .
<b>Llama 3.2 3B Instruct</b>	<b>Baseline:</b> According to the standard drill size chart, Drill Size 82 corresponds to a diameter of <b>0.820 inches</b> .	<b>Baseline:</b> Drill Size 84 refers to a drill bit with a diameter of 84 mm. To convert this to a decimal equivalent, we simply need to express the diameter as a decimal number. In this case, the decimal equivalent of Drill Size 84 is <b>84.00 mm</b> .	<b>Baseline:</b> The decimal equivalent of Drill Size 89 is <b>0.89 inches, or 0.0089 inches</b> .
	<b>Knowledge Graph:</b> To find the decimal equivalent of drill size 82, we need to look at the information provided. Drill size 82 corresponds to a decimal equivalent of <b>.0125 inches</b> .	<b>Knowledge Graph:</b> The decimal equivalent of Drill Size 84 is <b>.0115 inches</b> .	<b>Knowledge Graph:</b> The decimal equivalent of Drill Size 89 is <b>.0091 inches</b> .
<b>Qwen 2.5 7B Instruct</b>	<b>Baseline:</b> “Drill Size 82” likely corresponds to a drill bit that measures 82 thousandths of an inch in diameter. This translates to <b>0.82 inches</b> when expressed as a decimal.	<b>Baseline:</b> The decimal equivalent of Drill Size 84 typically refers to a drill bit size in the United States customary system... Rounded to four decimal places, the decimal equivalent of Drill Size 84 is approximately <b>0.0119 inches</b> .	<b>Baseline:</b> The decimal equivalent of Drill Size 89, which typically refers to a drill bit size in the United States customary system, is <b>1/2 inch</b> .
	<b>Knowledge Graph:</b> The decimal equivalent of Drill Size 82 is <b>.0125 inches</b> .	<b>Knowledge Graph:</b> The decimal equivalent of Drill Size 84 is <b>.0115 inches</b> .	<b>Knowledge Graph:</b> The decimal equivalent of Drill Size 89 is <b>.0091 inches</b> .

for models such as GPT-4o, GPT-4o mini, Gemini 2.0 Flash, Gemini 2.0 Flash-Lite, Llama 3.1 8B Instruct, and Llama 3.2 3B Instruct having an increase of upwards of  $3.1\times$  for ROUGE-1,  $2.8\times$  for ROUGE-2, and  $2.9\times$  for ROUGE-L. These substantial multiplicative gains clearly indicate that the integration of graph-derived information greatly enhances the models' ability to capture and reproduce relevant ground information thereby significantly improving the quality and fidelity of generated responses in content specific questions. Figure 6b) showcases an even more pronounced improvement for open ended generation when testing on the quantitative machining specific queries. For semantic similarity, there is an improvement up to  $1.2\times$  across all models with the exception of Gemini 2.0 Flash-Lite. For ROUGE metrics, all models exhibit an increase in scores compared to baseline with up to  $8.5\times$ ,  $8.8\times$ , and  $8.1\times$  for ROUGE-1, ROUGE-2, and ROUGE-L, respectively. These results provide evidence that more quantitative questions, which demand precise numerical accuracy and technical detail, benefit particularly from the enriched context. While semantic similarity indicates that the basic understanding of quantitative content is already strong, the substantial improvements in ROUGE scores reveal that the additional context helps the models generate far more accurate and detailed quantitative responses.

Table 4: Open ended responses from closed source models Gemini 2.0 Flash, Gemini 2.0 Flash-Lite, GPT-4o, and GPT-4o-mini regarding drill sizes for CNC machining. Grounded quantitative answers are highlighted in yellow, baseline values in red, and knowledge graph augmented values in green.

Model / Question	What is the decimal equivalent of Drill Size 82?	What is the decimal equivalent of Drill Size 84?	What is the decimal equivalent of Drill Size 89?
<b>Grounded Answer</b>	The decimal equivalent of Drill Size 82 is <b>.0125 inches</b> .	The decimal equivalent of Drill Size 84 is <b>.0115 inches</b> .	The decimal equivalent of Drill Size 89 is <b>.0091 inches</b> .
<b>Gemini 2.0 Flash</b>	<b>Baseline:</b> Drill sizes are often referred to by a number, letter, or fractional inch. Drill size 82 is a numbered drill size. Numbered drill sizes range from size 80, which is <b>.0135 inches</b> , to size 1, which is <b>.228 inches</b> . The size increases as the number decreases. Drill size 82 is equal to <b>.0156 inches</b> .	<b>Baseline:</b> The decimal equivalent of Drill Size 84 is <b>.0095 inches</b> .	<b>Baseline:</b> The decimal equivalent of Drill Size 89 is <b>.2010 inches</b> .
	<b>Knowledge Graph:</b> The decimal equivalent of Drill Size 82 is <b>.0125 inches</b> .	<b>Knowledge Graph:</b> The decimal equivalent of Drill Size 84 is <b>.0115 inches</b> .	<b>Knowledge Graph:</b> The decimal equivalent of Drill Size 89 is <b>.0091 inches</b> .
<b>Gemini 2.0 Flash-Lite</b>	<b>Baseline:</b> The decimal equivalent of drill size 82 is approximately <b>.0400 inches</b> .	<b>Baseline:</b> The decimal equivalent of Drill Size 84 is <b>.0400 inches</b> .	<b>Baseline:</b> The decimal equivalent of drill size 89 is approximately <b>.0400 inches</b> .
	<b>Knowledge Graph:</b> The decimal equivalent of drill size 82 is <b>.0125 inches</b> .	<b>Knowledge Graph:</b> The decimal equivalent of drill size 84 is <b>.0115 inches</b> .	<b>Knowledge Graph:</b> The decimal equivalent of drill size 89 is <b>.0091 inches</b> .
<b>GPT-4o</b>	<b>Baseline:</b> For drill size 82, which is part of the numbered drill bit system, the decimal equivalent is <b>.0120 inches</b> .	<b>Baseline:</b> In the standard number drill size system, the smallest size is 80, which corresponds to a diameter of <b>.0135 inches</b> . If you are referring to a drill size outside of this standard range, it might be a typographical error or a non-standard size specific to a particular manufacturer or application.	<b>Baseline:</b> In the context of standard numbered drill bits, Drill Size 89 is not typically recognized. The standard numbered drill bit sizes go from size 80, which has a decimal equivalent of <b>.0135 inches</b> .
	<b>Knowledge Graph:</b> The decimal equivalent of Drill Size 82 is <b>.0125 inches</b> .	<b>Knowledge Graph:</b> The decimal equivalent of Drill Size 84 is <b>.0115 inches</b> .	<b>Knowledge Graph:</b> The decimal equivalent of Drill Size 89 is <b>.0091 inches</b> .
<b>GPT-4o-mini</b>	<b>Baseline:</b> Drill Size 82 has a diameter of approximately <b>.4130 inches</b> .	<b>Baseline:</b> Drill size 84 corresponds to a diameter of <b>.0860 inches</b> .	<b>Baseline:</b> Drill size 89 corresponds to a diameter of approximately <b>.0135 inches</b> .
	<b>Knowledge Graph:</b> The decimal equivalent of Drill Size 82 is <b>.0125 inches</b> .	<b>Knowledge Graph:</b> The decimal equivalent of Drill Size 84 is <b>.0115 inches</b> .	<b>Knowledge Graph:</b> The decimal equivalent of Drill Size 89 is <b>.0091 inches</b> .

In addition to using quantitative metrics such as semantic similarity and ROUGE scores to evaluate each model's responses, we also include the full answers generated for the quantitative questions. This approach not only highlights the differences between closed source and open source models but also offers deeper insight into how each model integrates the retrieved graph knowledge into its responses. Table 3 and Table 4 show the baseline and knowledge graph-augmented outputs for drill size questions, which are focused on accurately determining the correct drill size

values for open and closed source models, respectively. Correct quantitative answers for each question are highlighted in yellow, baseline LLM values in red, and knowledge graph–augmented values in green for each table. As shown, the use of the knowledge graph consistently improves all open ended outputs of the LLMs’ answers by correcting the hallucinations present in the baseline answers. For example, in the open source group, models such as Llama 3.2 3B and Llama 3.1 8B provided baseline answers that deviated significantly from the grounded values (e.g., reporting 0.820 inches or 0.082 inches for drill size 82), whereas, when augmented with the knowledge graph, these models uniformly converged to the correct decimal equivalent of 0.0125 inches. Similarly, Qwen 2.5 7B Instruct confuses the drill size number for its equivalence in inches but this was effectively remedied with information from the graph. The closed sourced models exhibited a similar trend with Gemini 2.0-flash, Gemini 2.0-flash-lite, GPT-4o, and GPT-4o-mini all showing incorrect quantitative answers. These answers varied from nearly correct values, such as GPT-4o’s output of 0.012 inches for drill size 82, to very inaccurate ones, including GPT-4o-mini’s output of 0.4130 inches for the same drill size. Comparing across open and closed source models, the smaller open sourced models can tend to confuse the drill size number with its corresponding decimal value, effectively mistaking the numeric label for the measurement itself. In contrast, the closed sourced models often produce entirely incorrect values. In some cases, such as with GPT-4o, the model claims that the specified drill size does not exist. These results emphasize the value of incorporating a knowledge graph in open ended scenarios where choices are not given to correctly convert drill sizing information.

Table 5: Open ended responses from open source models Llama 3.2 3B Instruct, Llama 3.1 8B Instruct, and Qwen 2.5 7B Instruct regarding parameter selection for CNC machining operation. Grounded quantitative answers are highlighted in **yellow**, baseline values in **red**, and knowledge graph augmented values in **green**.

Model / Question	For drilling operations on steel (4140), what cutting speed (SFM) is recommended?	For tapping operations on aluminum, what cutting speed (SFM) is recommended?	For milling operations on brass with a tool diameter between 0.5 and 1.0 inches, what chip load is recommended?
<b>Grounded Answer</b>	For drilling operations on steel (4140), a cutting speed of <b>90 SFM</b> is recommended.	For tapping operations on aluminum, a cutting speed of <b>100 SFM</b> is recommended.	For milling operations on brass with a tool diameter between 0.5 and 1.0 inches, a chip load of <b>0.004 inches per revolution</b> is recommended.
<b>Llama 3.1 8B Instruct</b>	<b>Baseline:</b> The recommended cutting speed for drilling operations on 4140 steel ranges from <b>50 to 300 SFM</b> , depending on the drill bit type and desired finish.	<b>Baseline:</b> In general, the recommended cutting speed for tapping aluminum ranges from <b>50 to 150 SFM</b> , depending on the specific conditions.	<b>Baseline:</b> A commonly cited range for chip loads on brass is between <b>0.002 and 0.008 inches per tooth (ipt)</b> .
	<b>Knowledge Graph:</b> For drilling operations on steel (4140), the recommended cutting speed is <b>90 sfm</b> .	<b>Knowledge Graph:</b> For tapping operations on aluminum, the recommended cutting speed is <b>100 sfm</b> .	<b>Knowledge Graph:</b> For milling operations on brass with a tool diameter between 0.5 and 1.0 inches, the recommended chip load is <b>0.004 ipr</b> .
<b>Llama 3.2 3B Instruct</b>	<b>Baseline:</b> For drilling operations on 4140 steel, the recommended cutting speeds typically range from <b>60 to 120 SFM</b> , with the optimal speed depending on the specific drilling conditions.	<b>Baseline:</b> A cutting speed of <b>30 to 50 FPM</b> is commonly recommended for tapping operations on aluminum alloys, but the optimal cutting speed may vary depending on the application, tooling, and machining conditions.	<b>Baseline:</b> For milling operations on brass, a general rule of thumb is to use a chip load that is between <b>0.002 and 0.005 inches</b> for a tool diameter between 0.5 and 1.0 inches.
	<b>Knowledge Graph:</b> For drilling operations on steel (4140), a cutting speed of <b>90 SFM</b> is recommended.	<b>Knowledge Graph:</b> For tapping operations on aluminum, a recommended cutting speed is <b>100 sfm</b> .	<b>Knowledge Graph:</b> Based on the provided context, for milling operations on brass with a tool diameter between 0.5 and 1.0 inches, the recommended chip load is <b>0.004 inches per revolution (ipr)</b> .
<b>Qwen 2.5 7B Instruct</b>	<b>Baseline:</b> The recommended SFM for drilling 4140 steel generally ranges between <b>25 and 35 feet per minute (FPM)</b> .	<b>Baseline:</b> For tapping operations on aluminum, a recommended cutting speed (Surface Footage per Minute, SFM) typically ranges between <b>100 to 200 SFM</b> .	<b>Baseline:</b> A common starting point for chip loads in this range is approximately <b>0.003 to 0.005 inches per tooth (ipm)</b> .
	<b>Knowledge Graph:</b> For drilling operations on steel (4140), a recommended cutting speed (SFM) is <b>90 sfm</b> .	<b>Knowledge Graph:</b> For tapping operations on aluminum, a recommended cutting speed (SFM) is <b>100 sfm</b> .	<b>Knowledge Graph:</b> For milling operations on brass with a tool diameter between 0.5 and 1.0 inches, a chip load of <b>0.004 inches per revolution (ipr)</b> is recommended.

Table 5 and Table 6 present open and closed sourced model responses, respectively, to quantitative questions regarding selection of correct process parameters, which is especially important for ensuring optimal machining performance, precise material removal, and enhanced tool longevity across varying operational conditions. Similar to the previously mentioned tables, correct quantitative answers for each question are highlighted in yellow, baseline LLM values in red,

Table 6: Open ended responses from closed source models Gemini 2.0 Flash, Gemini 2.0 Flash-Lite, GPT-4o, and GPT-4o-mini regarding parameter selection for CNC machining operation. Grounded quantitative answers are highlighted in yellow, baseline values in red, and knowledge graph augmented values in green.

Model / Question	For drilling operations on steel (4140), what cutting speed (SFM) is recommended?	For tapping operations on aluminum, what cutting speed (SFM) is recommended?	For milling operations on brass with a tool diameter between 0.5 and 1.0 inches, what chip load is recommended?
<b>Grounded Answer</b>	For drilling operations on steel (4140), a cutting speed of <b>90 SFM</b> is recommended.	For tapping operations on aluminum, a cutting speed of <b>100 SFM</b> is recommended.	For milling operations on brass with a tool diameter between 0.5 and 1.0 inches, a chip load of <b>0.004 inches per revolution</b> is recommended.
<b>Gemini 2.0-flash</b>	<b>Baseline:</b> For 4140 steel, which is a medium-carbon, low-alloy steel known for its high strength and good hardenability, typical cutting speeds for drilling operations generally range from <b>70 to 120 SFM</b> .  <b>Knowledge Graph:</b> For drilling operations on steel (4140), a cutting speed of <b>90 SFM</b> is recommended.	<b>Baseline:</b> For general-purpose tapping of aluminum alloys like 6061-T6, a common starting point is between <b>30 to 60 SFM</b> . However, softer alloys like 1100 or 3003 can often be tapped at higher speeds, potentially up to <b>75 SFM</b> or more.  <b>Knowledge Graph:</b> For tapping operations on aluminum, a cutting speed of <b>100 SFM</b> is recommended.	<b>Baseline:</b> The recommended chip load for milling brass with a tool diameter between 0.5 and 1.0 inches typically falls within the range of <b>0.001 to 0.004 inches</b> per tooth.  <b>Knowledge Graph:</b> For milling operations on brass with a tool diameter between 0.5 and 1.0 inches, a chip load of <b>0.004 inches per revolution (ipr)</b> is recommended.
<b>Gemini 2.0-flash-lite</b>	<b>Baseline:</b> For high-speed steel (HSS) drills, a starting point for the cutting speed when drilling 4140 steel is typically in the range of <b>40 to 60 surface feet per minute</b> (SFM) for the steel in a hardened state.  <b>Knowledge Graph:</b> For drilling operations on steel (4140), the recommended cutting speed is <b>90 SFM</b> .	<b>Baseline:</b> For most aluminum alloys, a starting point for the cutting speed is typically in the range of <b>50 to 200 surface feet per minute (SFM)</b> .  <b>Knowledge Graph:</b> For tapping operations on aluminum, a cutting speed of <b>100 SFM (Surface Feet per Minute)</b> is recommended.	<b>Baseline:</b> For milling brass with a tool diameter between 0.5 and 1.0 inches, a chip load of <b>0.003 to 0.010 IPT</b> is a good starting point.  <b>Knowledge Graph:</b> For milling operations on brass with a tool diameter between 0.5 and 1.0 inches, a chip load of <b>0.004 inches per revolution (ipr)</b> is recommended.
<b>GPT-4o</b>	<b>Baseline:</b> For annealed AISI 4140 steel, which typically has a hardness of around 150-200 Brinell, the recommended cutting speed is usually in the range of <b>70 to 100 surface feet per minute (SFM)</b> when using high-speed steel (HSS) drill bits.  <b>Knowledge Graph:</b> For drilling operations on steel (4140), a cutting speed of <b>90 SFM</b> is recommended.	<b>Baseline:</b> The recommended cutting speed for tapping aluminum typically ranges from <b>75 to 150 surface feet per minute (SFM)</b> .  <b>Knowledge Graph:</b> For tapping operations on aluminum, a cutting speed of <b>100 surface feet per minute (SFM)</b> is recommended.	<b>Baseline:</b> For general milling operations on brass with a tool diameter between 0.5 and 1.0 inches, a typical chip load recommendation would range from approximately <b>0.002 to 0.010 inches per tooth (IPT)</b> .  <b>Knowledge Graph:</b> For milling operations on brass with a tool diameter between 0.5 and 1.0 inches, the recommended chip load is <b>0.004 inches per revolution (ipr)</b> .
<b>GPT-4o-mini</b>	<b>Baseline:</b> Typically, for drilling operations on 4140 steel, a cutting speed in the range of <b>60 to 100 SFM</b> is often recommended when using high-speed steel (HSS) drill bits.  <b>Knowledge Graph:</b> For drilling operations on steel (4140), the recommended cutting speed is <b>90 surface feet per minute (SFM)</b> .	<b>Baseline:</b> For most aluminum alloys, a common recommendation for tapping speeds is in the range of <b>100 to 300 SFM</b> .  <b>Knowledge Graph:</b> For tapping operations on aluminum, the recommended cutting speed is <b>100 surface feet per minute (SFM)</b> .	<b>Baseline:</b> For brass, which is a relatively soft and ductile material, the recommended chip load generally falls within the range of <b>0.002 to 0.008 inches per tooth</b> .  <b>Knowledge Graph:</b> For milling operations on brass with a tool diameter between 0.5 and 1.0 inches, the recommended chip load is <b>0.004 inches per revolution (ipr)</b> .

and knowledge graph-augmented values in green for each table. For the open sourced models, the baseline outputs typically provide a broad range of values such as Llama 3.2 3B Instruct suggesting a range from 60 to 120 surface feet per minute (SFM), and Llama 3.1 8B Instruct suggesting a range of 50 to 300 SFM. These wide ranges reflect uncertainty and a reliance on generic guidelines which do not allow a user to determine the optimal process parameters specified in technical documents. However, when the LLMs are augmented with knowledge graph information, these models converge on a precise recommendation of 90 SFM for drilling 4140 steel, 100 SFM for tapping aluminum, and a chip load of 0.004 inches per revolution (IPR) for milling brass, correctly aligning with the grounded answers. This shift from broad ambiguous ranges to a single well-defined value highlights the corrective influence of domain-specific information. The closed source LLMs exhibit a similar trend, frequently providing parameter recommendations as

broad ranges that mirror inherent uncertainty and a lack of industry-specific technical standards. However, when augmented with the knowledge graph the models once again converge to the correctly recommended process parameters. Comparing across open and closed models reveals that all tend to recommend broad and imprecise parameter ranges for operations such as drilling, tapping, and milling. Regardless of model size or architecture, the baseline outputs consistently display wide ranges that fail to account for material properties and precise machining requirements. This suggests that the inherent complexity and scale of the models do not necessarily translate into accurate internal representations of domain-specific machining standards. In summary, without integrating knowledge graph information the outputs remain too generic to guide optimal parameter selection demonstrating that the incorporation of structured, domain-specific data is essential for constraining outputs and calibrating process parameters to align with validated manufacturing standards.

## 6 Conclusion

This paper introduces ARKNESS (Augmented Retrieval Knowledge Network Enhanced Search and Synthesis), a hybrid LLM-agnostic pipeline that couples a semantically rich automatically constructed knowledge graph with any large language model to provide answers that are grounded and numerically precise. ARKNESS provides traceable answers validated in technical information through paragraph level entity-relation extraction that converts heterogeneous document formats into triples while retaining their original contextual information. Then, by storing the knowledge graph in a relational database, we can leverage vector-embedding search alongside beam search to retrieve the most relevant triples and contextual information, thereby conditioning the LLM to deliver precise, technically validated answers to user queries. Experiments on multiple choice questions across content specific and machining specific queries showcases uplifts up to 25% in accuracy and 22.4% in F1-score. For opened ended questions, uplifts up to  $8.1 \times$  for ROGUE-L scores demonstrate that the framework improves the structural and semantic information of generated responses relative to the grounded answers. These advantages allow the elimination of hallucinations prevalent in LLMs thereby enhancing process planning and usefulness in CNC machining contexts. Future work will focus on extending ARKNESS to a heterogeneous multimodal knowledge graphs, further strengthening factual grounding to enable closed-loop decision support in advanced manufacturing environments.

## Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## Acknowledgments

This work was supported under Cooperative Agreement W56HZV-21-2-0001 with the US Army DEVCOM Ground Vehicle Systems Center (GVSC), through the Virtual Prototyping of Autonomy Enabled Ground Systems (VIPR-GS) program.

DISTRIBUTION STATEMENT A. Approved for public release; distribution is unlimited. OPSEC9768

## References

- [1] M. Wiessner, P. Blaser, S. Böhl, J. Mayr, W. Knapp, K. Wegener, Thermal test piece for 5-axis machine tools, *Precision Engineering* 52 (2018) 407–417.
- [2] D. Hoang, H. Errahmouni, H. Chen, S. Rachuri, N. Mannan, R. ElKharboutly, M. Imani, R. Chen, F. Imani, Hierarchical representation and interpretable learning for accelerated quality monitoring in machining process, *CIRP Journal of Manufacturing Science and Technology* 50 (2024) 198–212.
- [3] Z. Chen, D. Hoang, R. Chen, F. Imani, Distributed hyperdimensional computing for real-time data aggregation and interpretable quality monitoring in manufacturing, in: ASME International Mechanical Engineering Congress and Exposition, volume 88605, American Society of Mechanical Engineers, 2024, p. V002T03A092.
- [4] Z. Li, Y. Dai, Z. Sun, C. L. Guan, T. Lai, H. Xu, X. Zhou, A sub-micron precision machining and measurement method of long travel metal guideways, *Journal of Manufacturing Processes* 133 (2025) 947–956.
- [5] ForInsights Consultancy, Computer numerical control (cnc) machine market forecast 2030, <https://www.forinsightsconsultancy.com/reports/computer-numerical-control-cnc-machine-market/>, 2023. Accessed: 2025-05-12.

- [6] D. Hoang, N. Mannan, R. ElKharboutly, R. Chen, F. Imani, Edge cognitive data fusion: From in-situ sensing to quality characterization in hybrid manufacturing process, in: International Manufacturing Science and Engineering Conference, volume 87240, American Society of Mechanical Engineers, 2023, p. V002T06A029.
- [7] D. Hoang, H. Chen, M. Imani, R. Chen, F. Imani, Brief paper: Multi-task brain-inspired learning for interlinking machining dynamics with parts geometrical deviations, in: International Manufacturing Science and Engineering Conference, volume 88117, American Society of Mechanical Engineers, 2024, p. V002T05A012.
- [8] K. Spanaki, D. Dennehy, T. Papadopoulos, R. Dubey, Data-driven digital transformation in operations and supply chain management, 2025.
- [9] Siemens AG, The True Cost of Downtime 2022, Technical Report, Siemens, 2022. URL: <https://assets.new.siemens.com/siemens/assets/api/uuid:3d606495-dbe0-43e4-80b1-d04e27ada920/dics-b10153-00-7600truecostofdowntime2022-144.pdf>.
- [10] N. C. Nwasuka, U. Nwaiwu, Computer-based production planning, scheduling and control: a review, *Journal of Engineering Research* 12 (2024) 275–280.
- [11] M. Raza, Z. Jahangir, M. B. Riaz, M. J. Saeed, M. A. Sattar, Industrial applications of large language models, *Scientific Reports* 15 (2025) 13755.
- [12] K. Šket, D. Potočnik, M. Ficko, S. Klančnik, Enhancing g-code programming in cnc machining using chatgpt: A comparative study of gpt-3.5 and gpt-4.0, Available at SSRN 4940034 (2024).
- [13] J. Jeon, Y. Sim, H. Lee, C. Han, D. Yun, E. Kim, S. L. Nagendra, M. B. Jun, Y. Kim, S. W. Lee, et al., Chatcnc: Conversational machine monitoring via large language model and real-time data retrieval augmented generation, *Journal of Manufacturing Systems* 79 (2025) 504–514.
- [14] R. Rosati, F. Antonini, N. Muralikrishna, F. Tonetto, A. Mancini, Improving industrial question answering chatbots with domain-specific llms fine-tuning, in: 2024 20th IEEE/ASME International Conference on Mechatronic and Embedded Systems and Applications (MESA), IEEE, 2024, pp. 1–7.
- [15] S. Kanimozhi, Y. Sriker, et al., Explorative deployment of fine-tuned large language model for on-site computerized numeric control machine operator assistance, in: 2024 IEEE Silchar Subsection Conference (SILCON 2024), IEEE, 2024, pp. 1–6.
- [16] Y. Xiao, S. Zheng, J. Shi, X. Du, J. Hong, Knowledge graph-based manufacturing process planning: A state-of-the-art review, *Journal of Manufacturing Systems* 70 (2023) 417–435.
- [17] P. Wen, Y. Ma, R. Wang, Systematic knowledge modeling and extraction methods for manufacturing process planning based on knowledge graph, *Advanced Engineering Informatics* 58 (2023) 102172.
- [18] D. Hoang, D. Gorsich, M. Castanier, F. Imani, Vector-symbolic knowledge graphs for enhanced memorization and reasoning in digital manufacturing, Available at SSRN 5097516 (2025).
- [19] L. Wang, H. Cheng, R. Wang, X. Huang, Machining scheme selection of features based on process knowledge graph and improved cosine similarity matching, *Machines* 13 (2025) 188.
- [20] J. Guo, J. Wu, J. Bian, Q. He, Knowledge graph-based machining process route generation method, in: International Conference on Human-Computer Interaction, Springer, 2023, pp. 35–48.
- [21] C. Cai, Z. Jiang, H. Wu, J. Wang, J. Liu, L. Song, Research on knowledge graph-driven equipment fault diagnosis method for intelligent manufacturing, *The International Journal of Advanced Manufacturing Technology* 130 (2024) 4649–4662.
- [22] Y. Li, H. Zhao, H. Jiang, Y. Pan, Z. Liu, Z. Wu, P. Shu, J. Tian, T. Yang, S. Xu, et al., Large language models for manufacturing, arXiv preprint arXiv:2410.21418 (2024).
- [23] P. Wang, J. Karigiannis, R. X. Gao, Ontology-integrated tuning of large language model for intelligent maintenance, *CIRP annals* 73 (2024) 361–364.
- [24] C. Qiu, B. Li, H. Liu, S. He, C. Hao, A novel method for machine tool structure condition monitoring based on knowledge graph, *The International Journal of Advanced Manufacturing Technology* 120 (2022) 563–582.
- [25] C. Auer, M. Lysak, A. Nassar, M. Dolfi, N. Livathinos, P. Vagenas, C. B. Ramis, M. Omenetti, F. Lindlbauer, K. Dinkla, et al., Docling technical report, arXiv preprint arXiv:2408.09869 (2024).
- [26] A. Hurst, A. Lerer, A. P. Goucher, A. Perelman, A. Ramesh, A. Clark, A. Ostrow, A. Welihinda, A. Hayes, A. Radford, et al., Gpt-4o system card, arXiv preprint arXiv:2410.21276 (2024).
- [27] M. Soori, F. K. G. Jough, R. Dastres, B. Arezoo, A review in capabilities and challenges of 5-axis cnc milling machine tool operations, Preprint (2024).

- [28] Y. Ye, T. Hu, C. Zhang, W. Luo, Design and development of a cnc machining process knowledge base using cloud technology, *The International Journal of Advanced Manufacturing Technology* 94 (2018) 3413–3425.
- [29] F. N. Guo, Exploring the application of industrial robots in cnc machining, *Journal of Global Humanities and Social Sciences* 4 (2023) 231–235. URL: <http://ojs.bonfuturepress.com/index.php/GHSS/article/view/1499>. doi:doi:10.61360/BoniGHSS232014990503.
- [30] Autodesk, Inc., Fundamentals of CNC Machining: A Practical Guide for Beginners, United States, 2014. URL: [https://academy.titansofcnc.com/files/Fundamentals\\_of\\_CNC\\_Machining.pdf](https://academy.titansofcnc.com/files/Fundamentals_of_CNC_Machining.pdf), desk Copy. Document Number: 060711.
- [31] I. Nugrahanto, H. Gunawan, H.-y. Chen, Innovative approaches to sustainable cnc machining: A machine learning perspective on energy optimization (2023).
- [32] M. Ihsan, Y. Sumantri, Y. Irawan, Integration of taguchi and promethee for cnc milling machining parameter optimization on aa6061, *International Journal of Mechanical Engineering Technologies and Applications* 5 (2024) 96–107.
- [33] S. Zhang, J. Bai, Research on cnc programming and machining process based on cad/cam technology, *Applied Mathematics and Nonlinear Sciences* 9 (2024) 1–18. URL: <https://doi.org/10.2478/amns-2024-0516>. doi:doi:10.2478/amns-2024-0516.
- [34] H. Wu, X. Wang, X. Deng, H. Shen, X. Yao, Review on design research in cnc machine tools based on energy consumption, *Sustainability* 16 (2024) 847.
- [35] M. Soori, F. K. G. Jough, R. Dastres, B. Arezoo, Robotical automation in cnc machine tools: a review, *acta mechanica et automatica* 18 (2024).
- [36] Meta AI, Llama 3.2 model card, [https://github.com/meta-llama/llama-models/tree/main/models/llama3\\_2](https://github.com/meta-llama/llama-models/tree/main/models/llama3_2), 2024. Version 3.2, released 25 Sept 2024.
- [37] A. Grattafiori, A. Dubey, A. Jauhri, A. Pandey, A. Kadian, A. Al-Dahle, A. Letman, A. Mathur, A. Schelten, A. Vaughan, et al., The llama 3 herd of models, *arXiv preprint arXiv:2407.21783* (2024).
- [38] A. Yang, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, C. Li, D. Liu, F. Huang, H. Wei, et al., Qwen2. 5 technical report, *arXiv preprint arXiv:2412.15115* (2024).
- [39] Google DeepMind, Gemini 2.0 flash, 2025. URL: <https://deepmind.google/technologies/gemini/flash/>.
- [40] Google DeepMind, Gemini flash lite, 2025. URL: <https://deepmind.google/technologies/gemini/flash-lite/>.