# A Survey of AI for Materials Science: Foundation Models, LLM Agents, Datasets, and Tools

**Minh-Hao Van** [*]
Department of EECS
University of Arkansas
Fayetteville, AR
haovan@uark.edu

**Prateek Verma** [*]
Department of EECS
University of Arkansas
Fayetteville, AR
prateek@uark.edu

**Chen Zhao**
Department of CS
Baylor University
Waco, TX
chen_zhao@baylor.edu

**Xintao Wu**
Department of EECS
University of Arkansas
Fayetteville, AR
xintaowu@uark.edu

## Abstract

Foundation models (FMs) are catalyzing a transformative shift in materials science (MatSci) by enabling scalable, general-purpose, and multimodal AI systems for scientific discovery. Unlike traditional machine learning models, which are typically narrow in scope and require task-specific engineering, FMs offer cross-domain generalization and exhibit emergent capabilities. Their versatility is especially well-suited to materials science, where research challenges span diverse data types and scales. This survey provides a comprehensive overview of foundation models, agentic systems, datasets, and computational tools supporting this growing field. We introduce a task-driven taxonomy encompassing six broad application areas: data extraction, interpretation and Q&A; atomistic simulation; property prediction; materials structure, design and discovery; process planning, discovery, and optimization; and multiscale modeling. We discuss recent advances in both unimodal and multimodal FMs, as well as emerging large language model (LLM) agents. Furthermore, we review standardized datasets, open-source tools, and autonomous experimental platforms that collectively fuel the development and integration of FMs into research workflows. We assess the early successes of foundation models and identify persistent limitations, including challenges in generalizability, interpretability, data imbalance, safety concerns, and limited multimodal fusion. Finally, we articulate future research directions centered on scalable pretraining, continual learning, data governance, and trustworthiness.

*Keywords* Materials Science · Foundation Models · LLM Agents · Datasets · Tools

## 1 Introduction

The field of materials science is entering a new era of data-driven discovery, accelerated by advances in artificial intelligence (AI) and machine learning (ML). Traditionally, computational materials science has relied heavily on first-principles simulations such as density functional theory (DFT), molecular dynamics (MD), and finite element methods to predict properties of materials, understand mechanisms, and guide experimental design. However, these methods are computationally intensive, grounded in approximations, and often constrained to small, well-characterized systems. In recent years, machine learning models trained on curated datasets, comprising both simulated and experimental results, have begun to supplement these traditional simulations, enabling faster property prediction and the emergence of generative design capabilities [1, 2, 3, 4, 5, 6]. Yet, most of these models remain task-specific, requiring dedicated architectures and training pipelines tailored to each property, material type, or data modality. Their generalization capacity, scalability, and cross-domain adaptability are thus limited.

Inspired by the transformative impact of foundation models (FMs) in natural language processing (NLP) (e.g., BERT [7], GPT [8, 9, 10], PaLM [11]) and computer vision (e.g., CLIP [12], DINO [13]), the materials science community is now exploring how similar large-scale and pretrained models might unlock new opportunities for research and innovation. Foundation models are typically defined as large, pretrained models trained on broad, diverse datasets and

---

[*]Both authors contributed equally to this research.

capable of generalizing across multiple downstream tasks with fine-tuning or prompt engineering. Their hallmark is the emergence of capabilities not explicitly programmed during training and the ability to transfer knowledge across domains, for example, from text to images or from property prediction to generative design, and aid in a variety of downstream tasks.

Materials foundation models aim to inherit these strengths while addressing the unique challenges of the physical sciences. First, materials data is inherently multimodal, comprising structures (atomic, crystalline, polymeric, and multiscale), textual descriptions, experimental data in the form of numbers, tables, and plots, images and spectra, experimental metadata, and simulated predictions or interpolations. Second, material properties emerge from the structure, assembly, and complex interaction of components at a variety of length scales ranging from subatomic, atomic, nanoscopic, microscopic, mesoscopic, and macroscopic. Third, many tasks require strict adherence to physical laws, such as energy conservation and symmetry constraints. Fourth, materials science suffers from limited labeled data; unlike NLP, it lacks billion-scale labeled corpora, relying instead on data that is costly to generate and often imbalanced.

Despite these challenges, recent advances illustrate the promise of this new paradigm. GNoME (Graph Networks for Materials Exploration) discovered over 2.2 million new stable materials by combining graph neural networks with active-learning-driven DFT validation [14]. MatterSim, a zero-shot machine-learned interatomic potential (MLIP), is trained on 17 million DFT-labeled structures and supports universal simulation across all elements and a wide range of temperatures and pressures [15]. MACE-MP-0, another universal MLIP, achieves state-of-the-art accuracy for periodic systems while preserving equivariant inductive biases [16]. Generative approaches such as MatterGen [17], DiffCSP++ [18], and CrystalFormer [19] enable conditional and multi-objective materials generation. Multimodal and cross-domain models like nach0 [20], MultiMat [21], and MatterChat [22] further demonstrate reasoning over complex combinations of structural, textual, and spectral data.

Other efforts aim to build generalist models that unify multiple domains and input types. ATLANTIC [23] explores cross-modal learning from literature, structures, and properties, while CrystaLLM [24] and GT4SD [25] provide frameworks for pretraining and multitask evaluation. Autonomous labs such as A-Lab integrate surrogate models and robotic synthesis to optimize experimental discovery [26]. Process-aware foundation models like Marcato FM [27] and applications to industrial-scale materials workflows [28] extend this paradigm to large-scale engineering systems and materials failure prediction.

Moreover, while initial FM research focused on crystalline, inorganic materials, there is growing recognition of the need to represent polymers, soft matter, disordered solids, and biological materials [29, 30]. These domains present challenges due to flexible, irregular, or long-range representations, and are motivating the design of new architectures, training data pipelines, and tokenization schemes. Early works such as AtomGPT [31] and MoL-MoE [32] are exploring this space.

Large language models (LLMs) are also being adapted for materials science. nach0 unifies natural and chemical language processing and performs tasks like molecule generation, retrosynthesis, and question answering [20]. Similarly, ChemDFM [33], LLaMat [34], and SciTune [35] are specialized LLMs trained on scientific literature and domain-specific data. These models enable tasks such as named entity recognition, synthesis extraction, literature summarization, and image-caption alignment.

LLM agents, which utilize LLMs as core reasoning components and interact with external environments, have been developed to support and automate tasks related to materials science. Recent studies have investigated the development of LLM-based agentic systems for materials science applications [36, 37, 38, 39, 40, 41]. HoneyComb [36] is designed to extend LLM capabilities in the materials science domain. LLMatDesign [37] is proposed to facilitate materials discovery by leveraging state-of-the-art LLMs. ChatMOF [38] presents an autonomous framework for predicting and generating metal-organic frameworks. MatAgent [39] is another LLM-based agentic system tailored for materials science, with a focus on property prediction, hypothesis generation, experimental data analysis, high-performance alloy and polymer discovery, data-driven experimentation, and literature review automation. MatPilot [40] focuses on literature search, scientific hypothesis generation, experimental scheme design, and autonomous experimental verification, with the goal of developing embodied AI capable of controlling physical robots.

Several toolkits and infrastructure platforms support this growing ecosystem. These include the Open MatSci ML Toolkit [42], designed for standardizing graph-based materials learning workflows, and FORGE [43], which provides scalable pretraining utilities across scientific domains. Combined, these efforts point to a future of deeply integrated, reusable, and generalizable AI systems for materials science.

In this paper, we present a comprehensive survey of foundation models in materials science. We categorize existing models by task, architecture type, and pretraining strategy; highlight multimodal and cross-domain models that bridge structure, text, and property spaces; summarize early successes and their implications for materials discovery pipelines; discuss open challenges such as data bias, long-range interaction modeling, and interpretability; and propose future
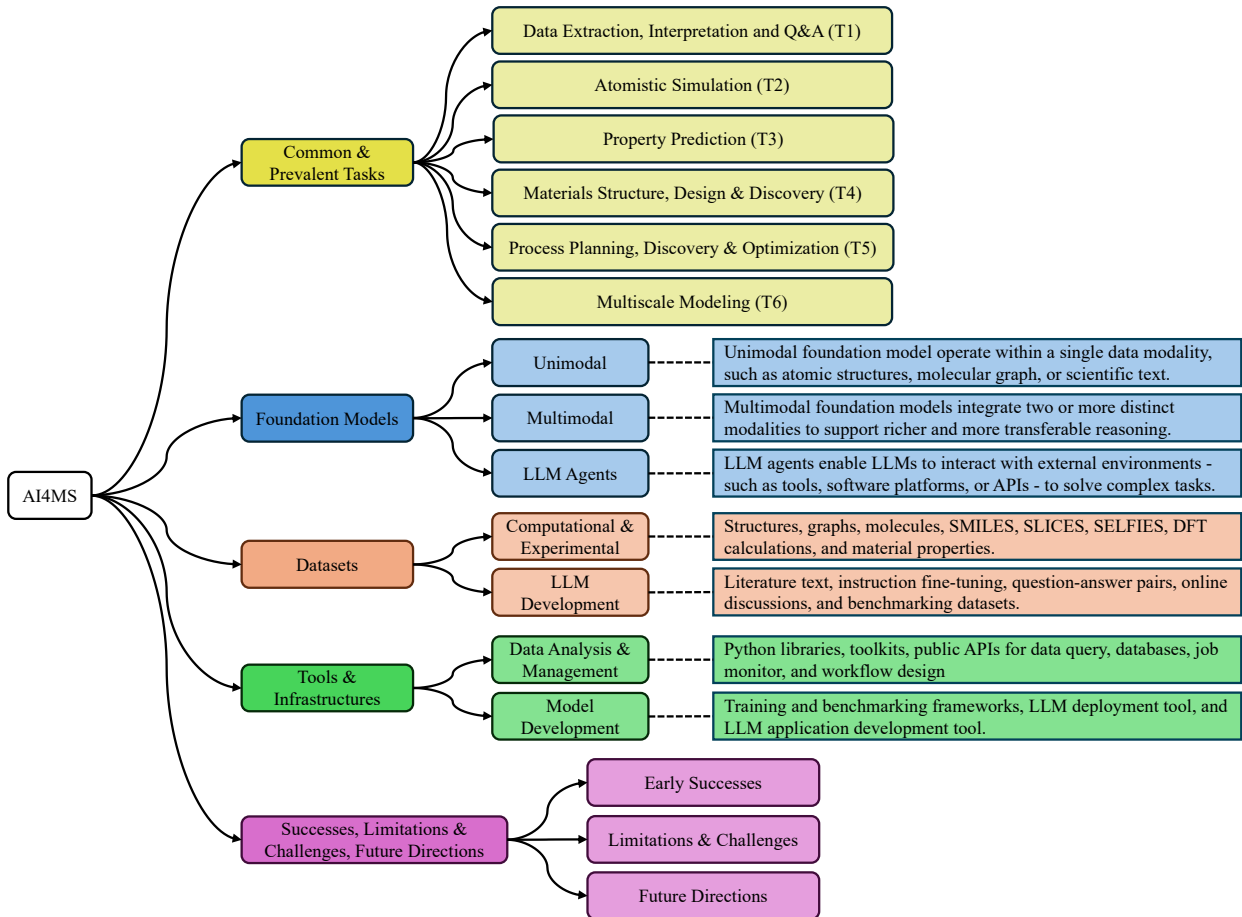
Figure 1: Overview of our survey of AI for materials science (AI4MS), highlighting common tasks, categories of foundation models, datasets, tools and infrastructures, as well as key discussions on early successes, current limitations, challenges, and future directions.

research directions toward scalable, multimodal, and human-AI collaborative systems. We structure this survey around a proposed taxonomy of foundation models in materials science, spanning six major application areas: data extraction, interpretation and Q&A; atomistic simulation; property prediction; materials structure, design and discovery; process planning, discovery, and optimization; and multiscale modeling. Figure 1 shows an overview of our survey of AI for materials science (AI4MS), including common materials science tasks, foundation models, datasets, useful tools and infrastructures for AI materials research, as well as key discussions of successes, limitations, challenges, and future directions. This work aims to serve both as a comprehensive reference and a forward-looking road map for researchers at the intersection of AI and materials science.

In compiling a comprehensive and in-depth survey of foundation models, datasets, and toolkits in materials science, we conducted a systematic literature search using Google Scholar[2]. Our search focused on research related to deep learning, foundation models, LLMs, and LLM-based agentic systems applied to materials science. We further extended the scope to include models in inorganic chemistry, with particular emphasis on crystal and atomic structures, which are highly relevant to materials research. We selectively incorporated high-quality research published in top-tier venues, as well as influential preprints from arXiv[3]. In addition to key papers, we include several open-source toolkits and resources that support the development of foundation models for materials research, selected based on our prior experience and their demonstrated utility in practical applications, as cited in the literature.

Prior to our work, several surveys have addressed the development and application of foundation models in materials science. Notably, Pyzer-Knapp et al. [1] focus on foundation models for materials discovery, categorizing them into

---

[2]Google Scholar, `https://scholar.google.com/`

[3]arXiv, `https://arxiv.org/`

four primary tasks: data extraction, property prediction, molecular generation, and synthesis prediction. In contrast, our survey encompasses a broader range of tasks relevant to materials science, including atomistic simulations and multiscale modeling. We systematically categorize models into unimodal foundation models, multimodal foundation models, and LLM agents, and further provide a comprehensive overview of available datasets and tools, offering readers additional resources to support and accelerate research in the field. Focusing specifically on crystalline materials, Wang et al. [2] present a survey of AI-accelerated approaches for crystal discovery, emphasizing four key tasks—property prediction, materials synthesis, characterization assistance, and acceleration of theoretical computations—as well as associated benchmarks, tools, and datasets. By comparison, our survey addresses a wider variety of material classes, including inorganic materials, organic compounds, polymers, and biomaterials, offering a more holistic view of AI applications across the materials science landscape. Another related survey by Han et al. [3] reviews recent advances in AI-driven inverse materials design, organizing models by material types or model architectures. However, this method of categorization may obscure a clear understanding of AI's progress in addressing practical, real-world challenges within specific materials science tasks. Expanding into the broader chemistry domain, Ramos et al. [4] discuss recent advances in LLMs and agentic AI for chemistry-related tasks. While there is considerable overlap between materials science and chemistry—particularly in tasks such as property prediction, synthesis planning, and information extraction—our survey remains firmly grounded in the core challenges and foundational models specific to materials science.

## 2   Common and Prevalent Tasks

In addition to categorizing foundation models by domain or architecture, it is useful to understand the broad functional tasks these models are designed to perform. Below, we organize key AI-driven tasks in materials science into six categories (T1–T6), describing where foundation models offer unique value and how these tasks span across material classes and length scales.

### 2.1   Data Extraction, Interpretation and Q&A (T1)

A significant portion of scientific knowledge in materials science is locked within unstructured data sources such as research papers, patents, lab notebooks, and experimental reports. The ability to read, interpret, and extract structured information from these documents is foundational to enabling data-centric discovery. Tasks in this area may include document classification, named entity recognition (e.g., identifying materials, properties, synthesis steps), synthesis route extraction, and scientific question answering. Foundation models, particularly large language models (LLMs) and multimodal Transformers, offer a generalizable framework across these varied tasks and data types with minimal retraining [1, 34, 29]. They can be tuned or prompted to extract synthesis protocols, material properties, designs and recommendations, or summarize literature trends at scale [26]. These efforts move beyond information retrieval and aim to construct structured, queryable knowledge graphs from raw scientific content.

### 2.2   Atomistic Simulation (T2)

Atomistic simulation tasks aim to replicate or accelerate quantum and molecular-scale simulations using AI. These include energy and force prediction, structure optimization, and molecular dynamics. Traditionally, each material or system required a custom force field or expensive ab initio calculations. Foundation models trained on millions of DFT-calibrated structures can now serve as general-purpose simulators across diverse chemistries and environments, offering near-DFT accuracy at a fraction of the computational cost [15, 16]. These models are especially effective for hard materials and periodic systems. However, they still face limitations in modeling long-range interactions, capturing rare events, and generalizing to non-equilibrium and disordered phases such as liquids, amorphous materials, or multi-component systems.

### 2.3   Property Prediction (T3)

Predicting material properties from composition or structure is one of the most mature and widely adopted tasks in materials science. This includes electronic, mechanical, thermal, optical, and chemical properties across a range of material classes. Traditional models often rely on domain-specific descriptors or task-specific neural networks trained on small datasets. Foundation models shift this paradigm by learning generalizable representations from large, diverse datasets, enabling transfer across tasks and domains [1, 29]. Once pretrained, these models can be fine-tuned or applied directly to new prediction tasks with minimal supervision. While most progress has centered on crystalline and small molecular systems, extending these models to disordered materials, porous frameworks, and polymers remains an ongoing challenge. Current models also focus primarily on equilibrium properties, while dynamic or temperature-dependent properties remain less explored [16].

## 2.4    Materials Structure, Design, and Discovery (T4)

Materials design encompasses tasks that involve generating or identifying new candidate materials with specified properties or performance metrics. This includes both molecular and extended solid-state systems, and can be conditioned on targets such as stability, conductivity, reactivity, or processability. Foundation models enable inverse design by learning structure-property relationships in reverse, allowing models to suggest new candidates given a desired property profile. Generative models, including diffusion models, graph-based Transformers, and LLM-driven molecular encoders, now support property-aware molecule and crystal generation, multi-objective optimization, and structure editing [17]. These approaches outperform traditional screening and optimization pipelines in navigating high-dimensional design spaces. However, synthesizability, dynamic stability, and real-world feasibility remain weak points, and polymer or disordered systems introduce additional complexity due to their flexible, long-range, and irregular representations.

## 2.5    Process Planning, Discovery, and Optimization (T5)

The successful realization of new materials requires not only design but also viable synthesis and processing pathways. Process-aware tasks include synthesis planning, reaction condition prediction, experimental design, and decision-making in closed-loop laboratories or production-line industries. Foundation models are increasingly used to extract and learn from synthesis protocols, suggest reaction pathways, and optimize experiment design in autonomous or semi-autonomous labs [26]. Tasks include recommending precursors, tuning process parameters, or guiding robotic experimentation. This category also encompasses broader workflows such as high-throughput screening, candidate ranking, uncertainty-aware exploration, and multi-property optimization. Foundation models provide scalable and generalizable interfaces for integrating these tasks into adaptive pipelines [34]. However, most applications remain constrained to well-studied inorganic materials. Incorporating real-world constraints like cost, scalability, toxicity, and manufacturability into these systems remains an open challenge, especially for polymers, biomaterials, composites, and other complex assemblies.

## 2.6    Multiscale Modeling (T6)

Materials performance often depends not only on atomistic configuration but also on structure and behavior at mesoscopic and macroscopic scales. Multiscale modeling tasks aim to integrate data across these length scales from atomic structure and microstructure to processing conditions, product-level properties, and performance degradation. Foundation models in this domain are still nascent but hold the potential to learn representations that capture temporal evolution, processing-structure-property relationships, and scale-bridging dynamics. Such models could complement finite element methods or serve as surrogates for time-dependent simulations. Current work focuses mostly on atomistic or crystal-level inputs—extending to grain boundaries, sintering behavior, fracture evolution, and large-scale failure remains underexplored [30]. Additionally, many technologically important materials such as polymers, gels, biological materials, and disordered solids do not conform to traditional crystalline representations. These materials present unique modeling challenges due to their flexible, irregular, and sequence-dependent structures. Foundation models could support new representations and generative strategies for these systems, including polymer sequence–property prediction, amorphous structure generation, and solvent-material interaction modeling. Developing domain-appropriate data, architectures, and training strategies for these systems is an important direction for future research.

# 3    Foundation Models in Materials Science

Foundation models (FMs) originate in natural language processing (NLP) and demonstrate an effective route to learning generalized representations through the mechanism of self-supervised training on a large corpus of text. In 2017, Vaswani et al. [44] first introduced the Transformer architecture based entirely on attention mechanisms. Attention operates by comparing each element in an input sequence to every other element. The strength of their association (called weights) captures relationships and contextual information within the sequence. The Transformer architecture is built around two key modules: the encoder and the decoder. The encoder consists of multiple stacked layers, and each layer adopts a self-attention mechanism. The input is tokenized from the model's vocabulary and and each token is represented as a vector that passes through the stacked layers. The decoder is composed of stacked layers, each containing a masked self-attention mechanism that ensures the model only attends to the current and previous tokens, preventing access to future tokens. Moreover, an encoder-decoder attention mechanism is adopted to align the decoder's output with relevant encoder inputs. The decoder can generate the sequence autoregressively, predicting each next token based on prior inputs. The encoder-decoder architecture is ideal for transforming sequences, such as translating from one language to another. The encoder itself can learn generalized representations and thus is ideal for tasks such as
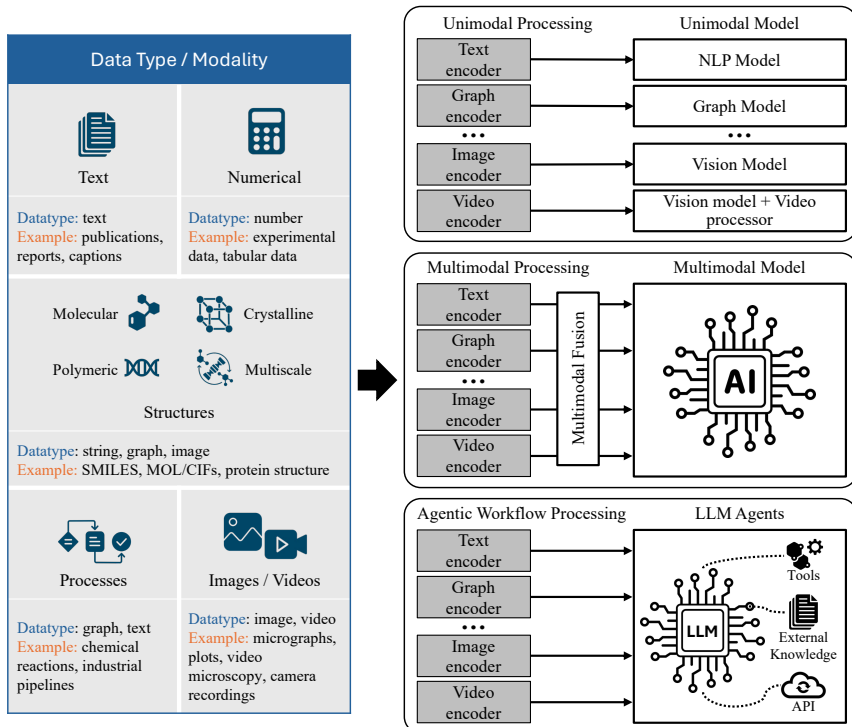
Figure 2: An illustrative example of the interplay of foundation models for materials science with data types and modalities.

property prediction whereas the decoder is ideal for tasks requiring sequence generation or completion such as inferring new outputs from input prompts.

Drawing from the success of the Transformer architecture, BERT as an encoder-only model, GPT as a decoder-only model, and BART as an encoder-decoder model further demonstrate impressive generalization across tasks such as summarization, translation, and question answering. BERT (Bidirectional Encoder Representations from Transformers) [7] was introduced in 2018 and utilizes only the encoder component. BERT's bidirectional Transformer is pretrained on the unlabeled text and the model processes the context both to the left and right of the word in question, thus enabling BERT to develop more comprehensive representations of input sequences, rather than mapping input sequences to output sequences. GPT (Generative Pretrained Transformer) [45] was introduced in 2018 as a decoder-only, left-to-right unidirectional language model to predict the next word in a sequence based on previous words, without an encoder. This decoder-based architecture forms the foundation for state-of-the-art large language models such as GPT-4 [10], Google Gemini [46], or Llama [47] models. BART (Bidirectional and Auto-Regressive Transformers) [48] was introduced in 2019 and is a sequence-to-sequence model consisting of a BERT-like bidirectional encoder and a GPT-like autoregressive decoder. BART is pretrained to reconstruct the original text from the corrupted ones (e.g., deleting tokens and shuffling sentences) with left-to-right autoregressive decoding as in GPT models. All these foundation models are trained on broad corpora using self-supervised objectives and then adapted to specific tasks with minimal fine-tuning. Their scale, in terms of parameters, data, and computing, enables emergent capabilities not explicitly programmed during training.

Beyond these static foundation models operating solely within their learned parameters and pretraining data, a new paradigm has emerged to support more autonomous and flexible problem-solving: LLM agents [49, 50, 51, 52, 53, 54, 55]. These systems enable LLMs to interact with external environments—such as tools, software platforms, external knowledge bases, or APIs—to solve complex tasks. In this paradigm, the LLM serves as a core reasoning agent: understanding the context, analyzing inputs, planning next steps, and taking actions by invoking external tools or resources. LLM agents thus offer a promising framework for interactive scientific discovery, capable of orchestrating design, prediction, simulation, synthesis planning, and information retrieval through prompting, planning, and tool use. These agentic capabilities open up new opportunities for materials science, where complex workflows often require sequential reasoning, hypothesis generation, interaction with simulations, or retrieval from literature text.

The concept of foundation models in materials science is relatively new and continues to evolve. A widely accepted definition describes FMs as large, pretrained models capable of learning transferable representations across diverse modalities and domains that can be reused across a range of downstream applications [29, 1]. These models differ from traditional machine learning pipelines, which are typically narrow, task-specific, and require extensive customization for each new task or data modality. In the context of materials science, FMs aim to unify representations across atomic structures, chemical formulas, natural language descriptions, spectroscopic data, and experimental metadata. The goals are to (i) enable transfer learning across tasks and materials systems, (ii) accelerate simulation and design workflows, and (iii) support generalist models for reasoning, generation, and prediction. Some FMs are relatively narrow (e.g., trained only for force prediction or molecular generation), while others aim to span multiple modalities and task types. These are sometimes referred to as "small" and "big" foundation models, respectively [29]. In this section, we first discuss unimodal foundation models tailored to individual data types, then explore multimodal approaches that integrate diverse inputs, and finally discuss emerging AI agentic frameworks that aim to assist with complex tasks in materials science. Figure 2 offers an illustrative example of how different types of foundation models process and learn from various data modalities in materials science research. To summarize key developments, Table 1 presents a comparative overview of notable foundation models, detailing their names, underlying architectures, supported modalities, datasets used, and associated tasks.

## 3.1 Unimodal Foundation Models

Unimodal foundation models operate within a single data modality, such as atomic structures, molecular graphs, or scientific text. These models are typically pretrained on large-scale datasets representative of their domain and often form the backbone of property prediction and simulation workflows.

### 3.1.1 Data Extraction, Interpretation and Q&A (T1)

A significant volume of materials information is contained in documents such as scientific publications, patents, and presentations. Extraction models with information retrieval capability have been extensively applied to identify materials from relevant documents and to link described properties with these materials.

For materials identification, **MatBERT** [70] is a BERT-based model trained for materials science, focusing on named entity recognition (NER) to extract and classify entities related to materials science into predefined labels. Being trained on a diverse set of NER datasets, MatBERT achieves impressive performance compared to baselines such as BERT. **MatSciBERT** [71] is a domain-specific, BERT-based model trained on a large corpus of peer-reviewed materials science publications and establishes state-of-the-art results on downstream tasks such as named entity recognition, relation classification, and abstract classification. **ChemDFM** [33] is a domain-specific language model trained on chemistry and materials literature. It supports tasks like synthesis information extraction, document classification, and literature-based reasoning.

For property extraction and association, **MaterialsBERT** [72] builds a general-purpose data extraction pipeline to automatically extract material properties from literature. MaterialsBERT is trained using 2.4 million materials science abstracts and obtains ~300,000 material property records that are made available at PolymerScholar[4]. Many tools for converting visual representations such as plots and charts into structured tabular data can help enhance the overall efficiency and accuracy of data extraction pipelines in materials. For example, **MolScribe** [73] is an image-to-graph generation algorithm to identify molecular structures from images in documents. It can predict atoms and bonds, along with their geometric layouts, for the molecular structure construction. **DePlot** [74] is a tool that translates the image of a plot or chart to a linearized table. The output of DePlot can then be directly used to prompt a pretrained LLM to extract necessary information.

A separate body of work examines the performance of existing or fine-tuned LLMs on materials science tasks or integrate them with other algorithms for materials science research. For example, Zaki et al. [75] present a dataset of 650 Q&A questions and evaluate the performance of Llama-2-70B, GPT-3.5, and GPT-4 models on solving these questions via zero-shot and chain of thought prompting. Similarly, Van Herck et al. [76] study the performance of fine-tuning three open-source LLMs (GPT-J-6B [77], Llama3.1-8B [78], and Mistral-7B [79]) for a range of different chemical questions, such as predicting properties of monomers from SMILES (Simplified Molecular Input Line Entry System) [80], and benchmark their performance against traditional binary classification models. **Uni-SMART** (Universal Science Multimodal Analysis and Research Transformer) [81] is designed for understanding multimodal scientific literature that contain molecular structures, chemical reactions, charts, and tables, in addition to textual content. The comparative evaluations demonstrate its superiority over several existing LLMs, such as GPT-4o, Gemini, and Claude [82]. **HoneyBee** [56] is an early attempt to build LLMs for materials science. A material-focused instruction

---

[4]PolymerScholar, `polymerscholar.org`

Table 1: Representative unimodal, multimodal, and agent-based foundation models in materials science. Each entry includes the model name, underlying architecture, supported modalities, datasets used, and associated AI tasks.

| Model Name | Architecture | Modality | Data Sources | Tasks |
|---|---|---|---|---|
| **Unimodal Foundation Models** | | | | |
| HoneyBee [56] | Llama-2 | Instruction fine-tuning (IFT) | MatSci-Instruct | T1 |
| ChemDFM [33] | Llama-2/3 | Literature text | PubChem, Q&A datasets, textbooks and papers | T1 |
| LLaMat [34] | Llama-2/3 | Literature text | Papers, RedPajama, MatSci community discourse | T1 |
| LLaMat-Chat [34] | LLaMat | Instruction fine-tuning (IFT) | OpenOrca, MathQA, MatSciNLP, MatBookQA, MaScQA | T1 |
| MACE-MP-0 [16] | Equivariant Graph Tensor Network | Atomic graphs | MPTrj | T2, T3 |
| MatterSim [15] | M3GNet + Graphormer | Atomic structures | MP, Alexandria, self-collected | T2, T3, T4, T5 |
| CGCNN [57] | GCNN | Crystal structures | MP | T3 |
| MEGNet [58] | GNN (distinct training for crystals or molecules) | Crystal structures, molecules | MP, QM9 | T3 |
| CatBERTa [59] | RoBERTa | Textual string of structures | OC20 | T3 |
| MoLXPT [60] | GPT-2 | Text, SMILES (text format) | PubMed, PubChem | T3 |
| GNoME [14] | GNN + Ensemble learning | Crystal graphs | MP | T3, T4, T5 |
| CDVAE [61] | VAE | Crystal structures | MP, Perov-5, Carbon-24 | T4 |
| DiffCSP [62] | Diffusion model | Crystal structures | MP, Perov-4, Carbon-24, MPTS-52 | T4 |
| MatterGen [17] | Diffusion model | Crystal structures | MP, Alexandria, ICSD | T4 |
| GP-MoLFormer [63] | Transformer | SMILES representations | ZINC, PubChem | T4 |
| CrystaLLM [24] | GPT-2 | CIFs | MP, OQMD, NOMAD | T4 |
| MatterGPT [64] | GPT-2 | SLICES representations | Alexandria | T4 |
| ChemFormer [65] | BART | SMILES | ZINC | T5 |
| FlowLLM [66] | Riemannian Flow Matching + LLMs | Textual strings of structures | MP | T4 |
| LLaMat-CIF [34] | LLaMat | CIF instruction fine-tuning | AMCSD, GNoME, MP | T4 |
| CSLLM [67] | Llama-2/3 | Crystal structures | ICSD, MP, COD, OQMD, JARVIS | T5 |
| **Multimodal Foundation Models** | | | | |
| LLM-Fusion [68] | GPT-2 | Text, SMILES, SELFIES, molecular fingerprints | QM9, ChEBI-20 | T3 |
| MoL-MoE [32] | Llama-3.2 | SMILES, SELFIES, molecular graphs | ZINC, ChEMBL, Moses | T3 |
| MatterChat [22] | Material encoder + Bridge model + LLM | Text, atomic structure | MPTrj, MP | T1, T3, T4, T5 |
| Text+Chem T5 [69] | T5 | Text, SMILES | Pistachio, CheBI-20 | T1, T4, T5 |
| nach0 [20] | T5 | Text, SMILES | Pubmed, USPTO, ZINC | T1, T4, T5 |
| SciTune [35] | LLaVA | Image, text | SciCap, ScienceQA | T1, T4 |
| MultiMat [21] | CLIP | Crystal structures, density of states, charge densities, text | MP, SNUMAT | T3, T4 |
| Marcato FM [27] | Encoder-Decoder + Llama-3 | Text, simulation grids | Self-collected simulations | T6 |
| **LLM Agents** | | | | |
| HoneyComb [36] | BM25 + Contriever + LLM | Literature text, web search | MatSciKB | T1 |
| LLMatDesign [37] | MatDeepLearn + TorchMD-Net + LLM | Structure, text | MP | T3, T4 |
| ChatMOF [38] | LLM as agent + LLM as evaluator + Tools | Structure, text | Collection of MOF databases | T1, T3, T4 |
| MatAgent [39] | LLM + ML models for materials science | Literature text, props | Self-collected data | T1, T3, T4 |
| MatPilot [40] | Knowledge retriever + physical workstations + LLM | Tabular data, text, graphs | N/A | T1, T5 |
| ChemCrow [41] | CoT reasoning + LLM | Literature search, web search | Online | T1, T4 |

fine-tuning framework, called **MatSci-Instruct**, is proposed to generate data for fine-tuning the Llama-based model. MatSci-Instruct utilizes ChatGPT as an instructor model to generate instruction-following data for materials science. Then, Claude serves as a verifier to assess the quality of generated data. Finally, s refinement-feedback loop approach is adopted to train HoneyBee with generated datasets from MatSci-Instruct. Other notable efforts include **LLaMat** and **LLaMat-Chat** [34], which are based on the Llama-2 and Llama-3 architectures. LLaMat comprises a family of models specialized in general materials science knowledge, distilled from research papers, textbooks, and online

forum discussions. LLaMat-Chat is an instruction fine-tuned variant of LLaMat, trained on both general and scientific question-answering datasets, thereby enabling more human-like, interactive dialogue capabilities.

### 3.1.2 Atomistic Simulation (T2)

In the domain of atomic-scale modeling, **MatterSim** [15] and **MACE-MP-0** [16] represent two large-scale efforts to build universal machine-learned interatomic potentials. Trained on tens of millions of DFT-labeled configurations, these models can simulate dynamics, phonon spectra, and phase stability across broad chemical spaces without system-specific retraining. While MatterSim adopts MEGNET [58] and Graphormer [83] to build the prediction pipeline, MACE-MP-0 is built on a state-of-the-art architecture, **MACE** [84], which utilizes equivariant and many-body message passing. **ANI** [85] utilizes Behler and Parrinello's symmetry functions [86] and neural network potentials to extract molecular representations, which are then utilized to predict molecular energy. Similarly, **AIMNet** [87] and **AIMNet2** [88] introduce atom-in-molecule networks to learn atomic representations and demonstrate their applications in predicting material properties.

### 3.1.3 Property Prediction (T3)

Identifying material properties is a complex task and classical ab initio physical simulation techniques, such as DFT, are computationally intensive. Deep learning techniques offer an efficient way to predict material properties by capturing the mapping relationship between materials data and properties from the collected datasets. Many of those models can be grouped into GNN-based methods [89] and Transformer-based methods. Specifically, these models are initially trained to encode the input features of materials (SMILES, SELFIES, crystal structures, etc.) into embedding vectors, which are then used to predict properties. It is important to note that SMILES and SELFIES (Self-Referencing Embedded Strings) [90] are 2D molecular representations that omit critical information about a molecule's 3D conformation. Widely used datasets such as ZINC [91, 92] and ChEMBL [93, 94] include billions of molecules but typically provide only 2D representations. In contrast, datasets of inorganic solids, such as crystalline materials, often include explicit 3D structural information, offering richer features for learning spatially dependent material properties.

Geometric GNNs are models designed to process graph data with geometric information such as spatial coordinates and angles. **SchNet** [95] is one of the earliest invariant GNN models designed for simulating quantum interactions in molecules using continuous-filter convolutional layers. This method directly models interactions between atoms utilizing distance information, thus achieving rotation invariant energy predictions. **CGCNN** [57] is designed for handling crystal structures and it can capture long-range interactions and global geometric structures. This model represents crystal structures as multi-edge graphs where nodes represent atoms and all their copies within the 3D space, and edges represent the connections between these atoms. CGCNN uses convolution and pooling layers to extract and learn both local and global features of the structure. CGCNN can extract the contributions from local chemical environments to global properties and thus is interpretable. **MEGNet** [58] is another GNN-based model and is trained on crystal structures from the Materials Project [96] or QM9 [97, 98] molecules, aiming to predict the formation energies, band gaps, or elastic moduli of crystals. Atomic attributes, bond attributes, and global state attributes serve as input features for MEGNet. In the pipeline, bond attributes are updated first, then atomic attributes, and lastly global state ones. A new graph representation is created from the MEGNet block. Additional set2set and dense layers are then added to predict the final output. We refer interested readers to a recent survey [99] that covers 80 GNNs used for property prediction.

Encoder-only Transformer architectures are primarily composed of an encoder, making them well-suited for property prediction that requires extracting meaningful information from input sequences such as SMILES. Encoder-only models based on the BERT architecture are predominantly used for property prediction as these models would ideally convert SMILES strings into a vector representation, which captures material properties. BERT-based models for property prediction include **ChemBERTa** [100], **CatBERTa** [59], **SolvBERT** [59], and **Mol-BERT** [101]. **ChemBERTa** [100] adapts RoBERTa [102] architecture to property prediction task using SMILES data from PubChem [103]. **CatBERTa** [59] is trained on a set of DFT calculations from Open Catalyst 2020 (OC20) [104] dataset. **SolvBERT** [59] is a model trained using an unsupervised scheme on solvation data. To effectively learn molecular representation, **SELFormer** [105] utilizes SELFIES data to train a Transformer-based RoBERTa model for extracting high-quality embeddings. Similarly, **Mol-BERT** focuses on learning molecular representation using SMILES from ZINC and ChEMBL datasets.

A line of research that utilizes GPT models have been introduced recently for predicting material property as Transformer-based architectures have shown remarkable capability in graph learning. **Matformer** [106] leverages the geometric distances between atoms from two adjacent unit cells to encode periodic patterns, thus enabling Matformer to encapsulate the lattice information and periodic patterns. **ComFormer** [107] converts both invariant graph representation and equivariant graph representation into embeddings via Transformers, enabling ComFormer to capture both local and global geometric information of different structures. **SMILES-GPT** [108] is a pretrained GPT-2-based language model

on a large SMILES corpus from PubChem. **MolXPT** [60] is introduced as a GPT-2-based language model using both text and SMILES from PubMed [109] and PubChem. **SPT** [110] takes SMILES as input to predict binary limiting activity coefficients, utilizing GPT-3 architecture. Another attempt to fine-tune GPT-3 on chemical and material data is introduced in [111].

### 3.1.4 Materials Structure, Design, and Discovery (T4)

This challenging task requires the development of advanced AI models to learn complex material structures and synthesize novel materials. Based on data representation, these methods can be grouped into two categories: geometric graph-based generation and string-based generation.

In the early stage, many GNN-based architectures have been proposed for this task. G-SchNet [112] is a pioneer in molecule generation. It incorporates the constraints of Euclidean space and the rotational invariances of the atom distribution as prior knowledge. It constructs an equivariant conditional probability distribution to determine the next atomic position by using the distance between the previously placed positions and the next atomic position as constraints. **CDVAE** [61] is proposed to generate stable crystalline materials from known materials, focusing on a generative approach using Diffusion [113] and VAE [114, 115]. The model includes three main components: (1) a periodic GNN encoder, (2) a property predictor, and (3) a periodic GNN decoder, which are optimized concurrently with stable materials. An evaluation of three tasks (reconstruction, generation, and property optimization) demonstrates the capability of CDVAE in generating materials given defined properties. Based on CDVAE, **Con-CDVEA** [116] generates crystals' latent variables according to given properties such as formation energy or band gap, and then yields the corresponding crystal structure by decoding the latent variables. **DiffCSP** [62] utilizes the fractional coordinate system to intrinsically represent crystals and model periodicity. By employing an equivariant graph neural network for the denoising process, DiffCSP conducts joint diffusion on lattices and fractional coordinates to capture the crystal geometry, thereby enhancing the modeling of the crystal geometry. DiffCSP separately and simultaneously adds noise and denoises on the fractional coordinates, atom types, and lattices. For the generation stage, random noises are sampled from Gaussian space for denoising as new materials. **DiffCSP++** [18] uses the diffusion model to generate new materials by incorporating the space group constraint. DiffCSP++ first encodes the lattice parameters as invariant vectors to ensure the lattice parameters satisfy the space group constraints. Then, DiffCSP++ generates the atom that satisfies the symmetry constrained by the space group.

**MatterGen** [17] is proposed to generate stable, diverse inorganic materials across the periodic table. MatterGen uses a diffusion model to produce crystalline structures by gradually refining atom types, coordinates, and the periodic lattice. Specifically, it uses a large dataset of stable material structures to train an equivariant score network and then fine-tunes the score network with a labeled dataset, where the property labels are encoded to steer the generation under specified property constraints. **GNoME** [14] is a graph neural network-based discovery engine and applies large ensembles and uncertainty quantification to predict material stability. It facilitates the discovery of new materials that extend beyond existing data distributions and has successfully discovered over 2 million new candidate crystals through active learning. The GNoME framework consists of two key modules: symmetry-aware partial substitutions combined with random structure search, and GNN-based modeling of material properties, and drive two independent materials discovery pipelines: structural pipeline and compositional pipeline. The former focuses on evaluating the stability of crystal frameworks without considering specific atomic types, filtering randomly generated structures using GNoME to retain potentially stable frameworks whereas the latter takes chemical formulas as input to GNoME, identifying stable chemical combinations to explore novel material compositions. DFT calculations are then performed to further validate their structural stability. Stable materials are added to the training set for subsequent iterations, creating an iterative active learning loop.

Crystalline materials are stored in standard text file formats known as CIF (Crystallographic Information File) or SLICES (Simplified Line-Input Crystal-Encoding System [117]). By treating CIFs or SLICES as plain string representations, many works explore the use of generative language models to generate crystals. These works [24, 118, 119] are different from models that use graph and graph-derived string representations as they treat materials as a sequence of discretized tokens and adopt Transformer architecture, thus utilizing the capability of next-token prediction given the input sequences for material generation. A majority of works for material structural generation are based on decoder-only GPT models. GPT employs positional encodings to maintain word order in its predictions. Its self-attention mechanism prevents tokens from attending to future tokens, ensuring each word prediction depends only on preceding words. **CrystaLLM** [24] is a decoder-only Transformer-based tool for crystal structure generation trained on an extensive corpus of the CIFs representing the structures of millions of inorganic solid-state materials. During training, the model is given a sequence of tokens from the corpus of CIFs, and is tasked with predicting the tokens that follow each of the given tokens. After training, the model can be used to generate new CIFs, conditioned on some starting sequence of tokens. For molecular design, **GP-MoLFormer** [63] enables property-conditioned molecule generation using Transformer architectures trained on large molecular datasets.

**CrystalFormer** [19] is a Transformer-based autoregressive model specifically designed to generate crystal materials that respect space group symmetries. It enriches the edge features' expressiveness by further incorporating angular information into the edge features and further introduces a graph construction method specifically designed for periodic invariance. **MatterGPT** [64] is another Transformer-based model trained for generating solid-state materials with given properties. Different from other language-based models which usually receive textual input, MatterGPT receives SLICES representation as input features and then generates output SLICES of expected materials. Flam-Shepherd et al. [119] demonstrate that language models trained directly on sequences derived directly from chemical file formats like XYZ files, CIFs, or Protein Data Bank (PDB) files can directly generate molecules, crystals, and protein binding sites in three dimensions. Gruver et al. [118] show that fine-tuned LLMs can generate the three-dimensional structure of stable crystals as text.

Being trained on a vast corpus of data, LLMs demonstrate impressive performance in multiple downstream tasks. **LLaMat-CIF** [34] is an instruction fine-tuned model with CIF data based on the pretrained LLaMat model. The goal is to enable the capability of LLMs in understanding and generating useful information with CIFs as inputs. **AtomGPT** [31] leverages GPT-2 and quantized Mistral models to learn the complex relationships between atomic structures and material properties from datasets such as JARVIS-DFT [120] and supports both property prediction and structure generation. While LLaMat-CIF and AtomGPT aim to train customized LLMs for materials science tasks, other works utilize strong pretrained LLMs to solve these tasks. **MatLLMSearch** [121] integrates pretrained LLMs with evolutionary search algorithms and supports crystal structure generation, crystal structure prediction, and multi-objective optimization of properties, all without fine-tuning. Beyond the direct use of LLMs to generate crystal representations, recent approaches have explored integrating LLMs with other generative techniques to enhance model performance. **FlowLLM** [66] combines LLMs with Riemann Flow Matching (RFM) to design novel crystalline materials. Specifically, FlowLLM first fine-tunes an LLM to learn a base distribution of meta-stable crystals in a text representation. These text representations are then converted into geometric graph representations. The RFM model takes samples from the LLM and iteratively refines the atom coordinates and lattice parameters to produce stable crystal structures. FlowLLM is trained on the widely used dataset of inorganic crystalline materials, derived from the Materials Project, focusing on a subset of compounds with up to 20 atoms known to be metastable.

### 3.1.5    Process Planning, Discovery, and Optimization (T5)

Synthesis planning aims for experimental realizations, i.e., predicting synthesizability and proposing the right precursors, pathways, and conditions to synthesize the targeted materials. **Molecular Transformer** [122] first applies a Transformer for synthesis prediction, in particular, translating reactants and reagents into the final product. This work studies the correlations between chemical motifs in reactants, reagents, and products in the USPTO dataset [123].

Advancements in synthesis prediction are mainly based on the BART encoder-decoder architecture [48]. **ChemFormer** [65] shows that models pretrained using only the encoder stack are limited for sequence-to-sequence tasks. ChemFormer is based on the BART encoder-decoder architecture and is trained with 100M SMILES from ZINC. The evaluations on datasets such as ChEMBL and ESOL [124] demonstrate state-of-the-art results in both sequence-to-sequence synthesis tasks and discriminative tasks.

Recent approaches have explored the use of LLMs for synthesis prediction. **MatChat** [125] demonstrates the effectiveness of LLMs in predicting the synthesizability of inorganic compounds and selecting suitable precursors. **CSLLM** [67] comprises three LLMs designed to predict material synthesizability, synthesis methods, and synthesis precursors. [67] employs Llama-7B fine-tuned via LoRA [126] and utilizes a proprietary Materials String representation to encode crystal structures. It further creates a synthesizability dataset containing 140,120 crystal structures. **SynAsk** [127] is an organic chemistry domain-specific LLM platform. SynAsk fine-tunes an LLM with domain-specific data and integrates it with a chain-of-thought approach. It supports functionalities such as molecular information retrieval, reaction performance prediction, and retrosynthesis prediction. **MatSci-LLM** [128] presents an LLM-based framework for generating experimental hypotheses of real-world materials discovery. By incorporating large-scale multimodal datasets for materials knowledge, MatSci-LLM framework can enable general-purpose LLMs to execute the materials discovery task via a six-step process: (1) materials query, (2) data retrieval, (3) materials design, (4) insilico evaluation, (5) experiment planning, and (6) experiment execution.

### 3.1.6    Multiscale Modeling (T6)

Materials simulations span length scales from atomic to macroscopic and time scales from femtoseconds to hours. Machine learning or deep learning based multi-scale modeling and fusion methods for materials research have been developed [129]. For example, **MuMMI** [130] is an ensemble machine learning approach for solving protein-membrane interactions that require multiscale modeling. This method connects three resolution scales: (1) coarsest scale (1000 nm), (2) coarse-grained scale (30 nm – 140K particles), and (3) all-atom scale (30 nm – 1.4M particles). CNN-based model

is utilized by [131] to achieve efficient multiscale modeling of heterogeneous materials, focusing on the prediction of homogenized macroscopic stress given a microstructure input image. However, there are no existing foundation model architectures specifically designed for multiscale modeling problems in materials science. This is partially because multiscale datasets are not yet well-established. Developing foundation models to establish connections and couplings across scales in materials remains an open challenge.

## 3.2 Multimodal Foundation Models

Multimodal foundation models integrate two or more distinct modalities including structure, text, spectra, and images to support richer and more transferable reasoning. These models are especially well-suited for capturing the interconnected nature of materials data, where a material may be simultaneously described by its atomic configuration, synthesis process, performance measurements, and visual or spectroscopic features.

**Text+Chem T5** [69] is a multi-domain, multi-task Transformer based language model and provides a unified representation between natural language and chemical representations. It is designed to support both mono-domain tasks and cross-domain tasks. The mono-domain tasks include molecule-to-molecule where the model predicts the outcome of a chemical reaction based on the starting chemicals or the starting chemicals that would be required to synthesize a given compound, and text-to-text where the model generates the action sequence based on a certain chemical reaction described in natural language. The cross-domain tasks include text-to-molecule where the model takes a textual description of a molecule as an input and generates its SMILES representation, and molecule-to-text where the model takes a molecule represented as SMILES and generates its human-readable textual description. The authors use multiple datasets including Pistachio [132] and CheBI-20 [133] in their evaluation. Another example is **nach0** [20], a multitask Transformer that unifies SMILES strings and natural language, enabling molecule generation, retrosynthesis, reaction prediction, and scientific question answering in a single model. **ATLANTIC** [23] extends this idea by aligning graph-based chemical representations with textual data, supporting synthesis reasoning and property prediction through interdisciplinary learning. **Regression Transformer** [134] abstracts regression as a conditional sequence modeling problem and bridges sequence regression and conditional sequence generation by using a nominal-scale training objective on combinations of numerical and textual tokens. It supports tasks of property prediction and conditional molecular design. Regression Transformer uses MoleculeNet [135] and TAPE (Tasks Assessing Protein Embeddings) [136] benchmarks in their evaluation.

**LLM-Fusion** [68] is a multimodal fusion model that leverages LLMs to integrate diverse representations, such as SMILES, SELFIES, text descriptions, and molecular fingerprints, for property prediction. Those diverse modalities are embedded and fused into a unified representation. In particular, the encoder for each modality can be of any architecture and can be frozen or fine-tuned. The encoded vectors via projection layers are stacked and enriched with positional encodings. The unified representation is fed at the input embeddings layer of the LLM, thus skipping the tokenization and positional encoding addition steps of transformer training. It uses the MoleculeNet-QM9 dataset and ChEBI-20 dataset. **MoL-MoE** [32] is introduced as a Multi-view Mixture-of-Experts framework to predict molecular properties by integrating latent spaces derived from SMILES, SELFIES, and molecular graphs. Mixture-of-Experts (MoE) has become essential for scaling large models by selectively activating sub-networks of experts through a gating network, thereby optimizing training efficiency. MoL-MoE utilizes SMI-TED (289M) foundation model [137] as the SMILES encoder, SELFIES-BART as the SELFIES encoder [138], and MHG-GNN [139] as the graph encoder. A set of nine distinct benchmark datasets sourced from MoleculeNet are used for both classification and regression tasks. **SciTune** [35] adopts a vision-language pretraining strategy, similar to CLIP [12], to align images (e.g., figures, microscopy) with captions and scientific questions. It supports multimodal tasks such as caption generation, figure interpretation, and visual question answering. Similarly, **MultiMat** [21] is a contrastive learning framework that learns embeddings across crystal structures, density of states (DOS), charge densities, and natural language labels.

**MatterChat** [22] is a versatile structure-aware multi-modal LLM that unifies material structural data and textual user queries and employs a bridging module to align a pretrained universal machine learning interatomic potential with a pretrained LLM. It supports text-based generation for tasks such as material property prediction, structural analysis, and descriptive language generation. MatterChat consists of three components: the Material Processing Branch for extracting atomic-level embeddings from material structural graphs, the Bridge Model for producing language model-compatible embeddings, and the Language Processing Branch for processing the user's text-based prompt into the language embeddings. Both language embeddings and query embeddings are fed into the LLM to produce the final text output. Finally, **Marcato FM** [27] is a cross-modal encoder-decoder model for material failure prediction, combining simulation outputs, structure, and grid-based fields for robust forecasting of fracture behavior.

Together, these models demonstrate the potential of foundation models to act as unified engines for materials understanding, integrating simulation, design, language, and experimentation within a single learning framework. However,

challenges remain in aligning modalities, addressing scale disparities, and ensuring physical fidelity across representation spaces.

### 3.3  LLM Agents

LLM agents can be designed to support and automate related tasks in materials science, including generating experimental plans, calling simulation tools, performing evaluation of outputs, and especially optimizing and repeating the experiments based on the outcome of previous experiment cycles.

**HoneyComb** [36] provides pretrained LLM capabilities to retrieve and analyze comprehensive knowledge in the materials science domain, enhancing the quality of generated content. HoneyComb consists of three main components: MatSciKB, a knowledge base aggregating diverse sources of materials-related information (e.g., arXiv articles, Wikipedia, datasets, textbooks); ToolHub, a suite of tools for retrieving up-to-date information; and Retriever, which extracts relevant knowledge from both MatSciKB and ToolHub using BM25 [140] and Contriever [141]. This agentic system can be integrated with different LLMs to enhance their domain-specific reasoning capabilities. **LLMatDesign** [37] employs a step-by-step, self-reflective design to accomplish the materials discovery task: (1) receiving human input about the chemical composition and the target property; (2) recommending the modification (i.e., exchange, addition, substitution, or removal); (3) employing machine learning tools such as MatDeepLearn [142] and TorchMD-Net [143] for property prediction; (4) evaluating the outcome; and (5) repeating the process with alternative modifications if the target property is not achieved. LLMatDesign translates human instructions into appropriate Materials Project API calls, supports material modifications, and evaluates outcomes using provided tools. Prompt engineering is used to guide pretrained LLMs in carrying out each step. This model randomly chooses ten starting materials from the Materials Project and focuses on designing materials with target properties of band gap and formation energy per atom, enabling a more efficient and autonomous approach to materials discovery. **ChatMOF** [38] is another agentic framework built to predict and generate metal-organic frameworks by leveraging pretrained LLMs, e.g., GPT-4, GPT-3.5-turbo, and GPT-3.5-turbo-16k, to extract key details from textual inputs and deliver appropriate responses. The system is comprised of three core components (i.e., an agent, a toolkit, and an evaluator) and it supports tasks of data retrieval, property prediction, and structure generation. **MatAgent** [39] aims to cover a sufficient number of tasks in materials science, such as property prediction, hypothesis generation, experimental data analysis, high-performance alloy and polymer discovery, data-driven experimentation, and literature review automation. Unlike other systems such as HoneyComb and LLMatDesign, MatAgent involves a human-in-the-loop (HITL) approach, allowing for human oversight at key stages to ensure the quality of generated content. The workflow consists of: (1) hypothesis generation and initial review with HITL; (2) central processing, including code generation, data visualization, report writing, and web search for prediction and analysis; (3) quality review to check for non-sense outcome with HITL and move back to previous step if it happens; (4) final human review; and (5) integration with external tools or databases. Similar to MatAgent, **MatPilot** [40] also employs human-in-the-loop approach, with a particular emphasis on literature search, scientific hypothesis generation, experimental scheme design, and autonomous experimental verification. MatPilot is designed with two primary modules: (1) a cognition module, which retrieves information from a knowledge base and generates experimental schemes under the guidance of human experts; and (2) an execution module, which autonomously performs experimental tasks based on the plan produced by the cognition module. Unlike previous agents designed to interact with tools or APIs, MatPilot moves toward interaction with physical robots, aiming to establish embodied AI capable of performing real-world experimental procedures. **ChemCrow** [41] aims to build an LLM agent for chemistry, covering several tasks including organic synthesis, drug discovery, and materials design implemented with the LangChain framework. This LLM-based agent utilizes a chain-of-thought reasoning loop to (1) create plans, (2) select appropriate tools, (3) take actions, and (4) analyze the outputs. Incorporating a set of 18 chemistry tools covering general tools, molecule tools, safety tools, chemical reaction tools, ChemCrow enables LLM with autonomous experimentation in chemistry tasks.

## 4  Datasets, Tools, and Infrastructures

### 4.1  Datasets

The effectiveness of foundation models in materials science hinges critically on the availability of high-quality, large-scale datasets that span diverse materials classes, properties, and modalities. Over the past few years, a growing number of datasets have emerged to support pretraining and evaluation, many of which are summarized in Table 2. These datasets vary widely in terms of modality - ranging from atomic structures, compositions, and energy landscapes to text, graphs, and synthesis protocols.

Table 2: Representative datasets used for training and evaluating foundation models in materials science (MatSci). The table covers both computational/experimental and LLM-oriented datasets, detailing their modalities, sample size, and relevance to key AI tasks (T1–T6).

| Dataset | Data Description | Modalities | # Samples | Tasks |
|---|---|---|---|---|
| **Computational and experimental datasets** | | | | |
| ICSD [144] | Inorganic materials | Crystal structures, props | ∼300k | T3, T4 |
| Materials Project (MP) [96] | Structure, property | Structured text, graphs | ∼200k | T2, T3, T5, T6 |
| OQMD [145, 146] | Inorganic crystals | Structure, DFT props | ∼1M | T2, T3 |
| NOMAD [147] | Hypothetical crystals | Crystal structures, props | ∼214k | T3, T4 |
| OMat24 [148] | Inorganic materials | Crystal structures, props | ∼118M | T3, T4 |
| SNUMAT [149] | Synthesized materials | Crystal structures, DFT props | ∼10k | T3, T4 |
| MPTrj [150] | Trajectory data | Atomic structures, props | ∼1.5M | T2, T3 |
| Open Catalyst 2020 (OC20) [104] | Catalysis structures, forces | Graphs, 3D coords | ∼1.3M | T2, T3, T5 |
| Alexandria [151] | 1D, 2D, 3D materials | Molecular structures, DFT calculations | ∼5M | T2, T3, T4 |
| QM9 [97, 98] | Molecules, DFT | SMILES, 3D structures | 134k | T2, T3 |
| Guacamol [152] | Benchmark for molecular design | Molecules | ∼1.5M | T4 |
| Moses [153] | Benchmark for molecular design | Molecules | ∼4.5M | T4 |
| ZINC [91, 92] | Molecular library | SMILES, conformers | >1B | T3, T4 |
| ChEMBL [93, 94] | Bioactive molecules | Molecules | ∼2.2M | T3, T4 |
| **LLM development datasets** | | | | |
| MatScholar [154, 155] | Literature text | Paper abstracts | >5M | T1 |
| MatSciKB [36] | Knowledge base | Text | ∼38k | T1 |
| PubChem [103] | Proteins, genes, chemical structures, literature text | SMILES, crystal structures, papers, patterns | >200M | T1, T3, T4 |
| MatbookQA [34] | Q&A from MatSci books | IFT, Q&A | ∼2k | T1 |
| MaScQA [75] | Q&A from engineering exams | IFT, Q&A | ∼1.5k | T1 |
| MatSci-Instruct [56] | MatSci IFT data | IFT, Q&A | ∼52k | T1 |
| MatSciNLP [156] | NLP benchmark | Text, Q&A | ∼170k | T1, T4, T5 |
| LLM4Mat-Bench [157] | LLM benchmark for property prediction | Crystal composition, CIFs, textual descriptions | ∼1.9M | T3 |
| LLM4Mol [158] | LLM benchmark for molecular prediction | Molecules, props | ∼40k | T4 |
| MACBENCH [159] | LLM multimodal benchmark for chemistry and MatSci | Q&A | 628 | T1 |

### 4.1.1 Computational and experimental datasets

In the last two decades, several computational and experimental datasets have been released to support materials science, which are later used to train deep learning and foundation models to further enhance materials exploration and analysis. In the early stage, **ICSD** [144] received its first record back in 1913 and was later available on the web in 2003. It is the largest database of inorganic crystal structures, comprising about 300,000 entries. Each entry includes detailed crystallographic information, such as unit cell parameters, space group, atomic coordinates, site occupation factors, molecular formulas, and molecular weights. After that, other common datasets for crystalline solids, material properties, and DFT calculations—such as **Materials Project (MP)** [96], **OQMD** [146, 145], **NOMAD** [147], **Open Material 2024 (OMat24)** [148], and **SNUMAT** [149]—have been released to mainly serve as training data for deep learning and foundation models. Specifically, MP provides an open database of over 200,000 materials through web interfaces and APIs. OQMD is another dataset with high-throughput DFT calculations of thermodynamic and structural properties of 1,317,811 materials. NOMAD serves as an open-access platform for managing, analyzing, and sharing materials science data, currently hosting approximately 214,000 structures. Recently, OMat24, released by FAIR, presents a large-scale database of inorganic materials accompanied by a set of pretrained models. It contains around 118 million structures labeled with properties such as total energy, atomic forces, and cell stress. To generate these crystal structures, the process begins with a random sampling of relaxed structures, followed by one of three refinement techniques: Rattled Boltzmann sampling, ab initio molecular dynamics (AIMD), or Rattled relaxation. SNUMAT provides API access to its database of around 10,000 synthesized materials and their DFT properties. Other datasets, such as **Open Catalyst 2020 (OC20)** [104] and **MPTrj** [150], have been designed with large-scale adsorption energies and structures on catalyst surfaces, interatomic potentials, and atomic trajectory. For exploring molecular property, **Alexandria** dataset [151] offers an open collection of approximately 5 million DFT-calculated molecular structures. The dataset is categorized by material type, including 3D, 2D, and 1D materials. **QM9** [97, 98] includes over 100,000 small organic molecules with DFT-calculated quantum chemical properties. To support molecular design, **Guacamol** [152] and **Moses** [153] are benchmarks for evaluating models trained for molecular generation task. Some chemical datasets can also serve as additional resources for retrieving specific information about materials. For example, **PubChem** [103] is the largest open database of chemical information including name, formula, structure, and other identifiers. The database hosts

around 121M chemical structures and hundreds of thousands of proteins and genes. **ZINC** [92, 91] is another chemical database of compounds for virtual screening, including over one billion of molecules. **ChEMBL** [93, 94] is a database for drug discovery, including bioactive molecules and drug-like properties.

Despite this progress, the landscape remains fragmented. Table 3 provides a cross-sectional view of widely used computational and experimental datasets in materials science with respect to the types of materials they support and the physical length scales they span. While most datasets are heavily skewed toward inorganic and atomically resolved systems, a few, particularly NOMAD and ChEMBL extend into organic, polymeric, and biomaterial domains. This uneven distribution highlights a critical gap in the availability of structured, large-scale data for polymers, composites, and mesoscale to macroscale phenomena. The dominance of atomistic simulation datasets has naturally shaped the focus of existing foundation models, many of which are optimized for crystalline inorganic compounds. As the field moves toward more generalizable and multimodal models, addressing this data imbalance is essential for expanding foundation model applicability across underexplored material classes and real-world scales. This challenge underscores the need for systematic data curation and benchmarking efforts that extend beyond atomic configurations, enabling models to reason across structure, property, processing, and scale. Moreover, few datasets provide truly multimodal, property-aligned examples—such as atomic structure paired with synthesis route and experimental spectra—thereby limiting the development of generalist models capable of reasoning across diverse data types. Many datasets also suffer from compositional, elemental, and phase-type biases, with overrepresentation of stable oxides and underrepresentation of disordered or metastable materials [15]. Ongoing efforts in data augmentation, transfer learning, and active learning aim to address these limitations by prioritizing data diversity, especially in low-resource domains like polymers, amorphous systems, and biomaterials [29, 30].

Table 3: Coverage of computational and experimental datasets across material types and length scales. Checkmarks indicate areas where the dataset or model is primarily applied or relevant.

| Dataset / Model | Material Type | | | | Length Scale | | | |
|---|---|---|---|---|---|---|---|---|
| | Inorganic | Organic | Polymers | Bio | Atomic $< 10^{-9}m$ | Nano $10^{-9} - 10^{-6}m$ | Meso $10^{-6} - 10^{-3}m$ | Macro $> 10^{-3}m$ |
| **Datasets** | | | | | | | | |
| ICSD [144] | ✓ | ✓ | | | ✓ | ✓ | | |
| Materials Project (MP) [96] | ✓ | | | | ✓ | ✓ | | |
| OQMD [145, 146] | ✓ | | | | ✓ | | | |
| NOMAD [147] | ✓ | ✓ | ✓ | ✓ | ✓ | ✓ | | |
| OMat24 [148] | ✓ | | | | ✓ | ✓ | | |
| SNUMAT [149] | ✓ | | | | ✓ | ✓ | | |
| MPTrj [150] | ✓ | | | | ✓ | ✓ | | |
| Open Catalyst 2020 (OC20) [104] | ✓ | | | | ✓ | ✓ | | |
| Alexandria [151] | ✓ | | | | ✓ | | | |
| QM9 [97, 98] | | ✓ | | | ✓ | | | |
| Guacamol [152] | ✓ | ✓ | | | ✓ | ✓ | | |
| Moses [153] | | ✓ | | | ✓ | ✓ | | |
| ZINC [91, 92] | | ✓ | | | ✓ | ✓ | | |
| ChEMBL [93, 94] | | ✓ | | ✓ | ✓ | ✓ | | |

### 4.1.2   LLM development datasets

With the rapid development of LLMs such as ChatGPT, Gemini, Llama, or Claude, more real-world AI-assisted systems have been built to bring useful domain-specific applications to users. The scientific domain, especially materials science, is an important goal in achieving artificial general intelligence. Developing LLMs requires several stages including training, instruction fine-tuning, and evaluation, which necessitate several large datasets for various tasks.

For training LLMs with general knowledge, **RedPajama** [160], **BookCorpus** [161], **Wikitext** [162] and **Common Crawl** [163] are common open-source datasets including billions of entries for LLM pretraining. Other datasets of literature papers and patterns are also used for training LLMs in scientific domains. **MatScholar** [154, 155] includes over five million paper abstracts in the materials science domain. **SciDocs** [164] offers a foundation for language model pretraining, focusing on scientific documents. **USPTO** [123] offers Open Data Portal (ODP) for accessing US patterns, which enables researchers to discover and extract related information to specific domains. In addition to chemical compounds, **PubChem** [103] also provides 53M patents and 42M research papers. All these literature-related datasets can sufficiently enable the training of LLMs for materials science. **SciCap** [165] is a multimodal dataset with scientific figure-caption pairs extracted from over 290,000 arXiv papers, which is suitable for training multimodal models like vision-language models. Comprehensive knowledge bases are also essential to operate LLM-based agentic systems as they require a sufficient amount of information to generate hypotheses, do planning, or provide accurate answers. To build LLM-based agents for materials science, **MatSciKB** [36], which is part of the HoneyComb work, includes

diverse sources of data related to materials science—such as arXiv papers, Wikipedia pages, textbooks, multiple-choice datasets, formulas, and GPT-generated entries—to support the task of knowledge retrieval, hypotheses generation, and experiment planning.

Instruction fine-tuning (IFT) is another critical step when developing LLM-based chat assistants, enhancing the capability of generating human-like conversations. Several datasets are introduced to support IFT including general Q&A like **OpenOrca** [166] and **WebInstructSub** [167], and science-related Q&A extracted from science exams or knowledge bases of different subjects, such as **MathQA** [168], **MatbookQA** [34], **MaScQA** [75], **ARC** [169], **PIQA** [170], **SciQ** [171] and **ScienceQA** [172]. **MatSci-Instruct** [56], which is proposed in HoneyBee work, is an attempt to build an instruction fine-tuning dataset specified for materials science with around 52k instructions. **MatSci community disclosure**[5] is a forum for discussions of topics related to materials science, which is also a helpful Q&A data source.

While several LLMs have been introduced for materials science, the cost of deploying an LLM on local machines is expensive due to the high cost of running GPU servers. Cloud-based pretrained LLMs, such as ChatGPT and Gemini, are still cost-efficient options for materials science researchers. Hence, the task of evaluating and benchmarking pretrained LLMs in the materials science domain has been catching the attention of the research community to assess the performance of pretrained models on downstream tasks related to materials science. To evaluate language models on materials science text, **MatSciNLP** [156] provides a benchmark comprising seven tasks, including Named Entity Recognition (NER), relation, sentence, and paragraph classification, event argument extraction, synthesis action retrieval, and slot filling. **LLM4Mat-Bench** [157] offers a benchmark specifically focused on evaluating LLM capabilities in property prediction, particularly for crystalline materials. It includes around 1.9M of structures formed by 10 public materials datasets with 45 distinct properties. Focusing on multimodalities, LLM4Mat-Bench covers several modalities, such as crystal compositions, CIFs, and corresponding textual descriptions. Several models are evaluated to show their performance in material property prediction tasks, such as CGCNN, MatBERT, LLM-Prop, Llama, Gemma, and Mistral. **LLM4Mol** [158] introduces a benchmark aimed at evaluating the ability of LLMs to perform molecular prediction tasks. The benchmark covers six molecule datasets in two important prediction tasks: classification and regression. Prompt engineering is utilized to instruct LLMs, such as GPT-family and Llama-family models, to predict outcomes using both zero-shot and few-shot prompting techniques. **MACBENCH** [159] is a multimodal benchmark comprising 628 questions designed to evaluate the multimodal reasoning capabilities of AI models in materials science and chemistry tasks. MACBENCH encompasses tasks of fundamental scientific understanding, data extraction from visual information, and practical knowledge. Its diverse visual inputs include laboratory images, band structures, crystal structures, and atomic force microscopy images paired with multiple-choice questions. This benchmark addresses the multimodal nature of materials science and captures the tacit knowledge and laboratory skills.

While the descent effort of introducing training and benchmarking datasets for materials science has been introduced, the need for comprehensive and multi-task datasets still remains. Furthermore, there is also a lack of multimodal datasets for training multimodal LLM.

## 4.2   Tools and Infrastructures

In this section, we review useful tools and infrastructures designed for the analysis and management of materials data, as well as for the development of machine learning models. A summary of key resources is provided in Table 4.

### 4.2.1   Materials data analysis and management tools

In the effort to enhance materials data analysis and management, several packages and libraries are developed to support materials science research. **Pymatgen** (Python Materials Genomics) [174] is an open-source Python library for processing, analyzing, and visualizing crystal structures, phase diagrams, and material properties in different formats (e.g., VASP, ABINIT, CIF, XYZ). A Python-based web app framework, named **Crystal Toolkit** [182], is introduced to provide an interactive interface for exploring materials science information. $M^2$**Hub** [175] is another Python-based toolkit that provides machine learning practice for materials science. It supports users with the whole pipeline creation, covering from data processing tasks (downloading and preprocessing) to model implementation and training. Several additional Python packages have been developed to design and manage end-to-end scientific workflows. These include **FireWorks** [176], **Custodian** [174], **Atomate** [183], and **Jobflow** [184], which facilitate tasks such as job writing, execution, automation, management, and result analysis. For the purpose of building a data pipeline and database, **Emmet** [185] and **Maggma** [177] are developed by Materials Project to support users with data queries, data transformations, and data storage. HoneyComb provides **ToolHub**, a unified interface consisting of several tools for accessing online sources such as Google search, Wikipedia search, or materials analysis tools. To build an API

---

[5]MatSci Com. Dis., `https://matsci.org/`

Table 4: Key tools and infrastructures supporting foundation models in materials science.

| Name | Primary Functionality | Supported Modalities | Supported Models/Applications | Access |
|---|---|---|---|---|
| **Materials data analysis and management tools** | | | | |
| Materials Project API [96] | Access to millions of computed materials entries for downstream use | Structures, formation energy, bandgap | Various via API | Open-source |
| OPTIMADE [173] | Access to millions of material structures and information | Structures, DFT calculations, references | Various via API, Python tools | Open-source |
| MPTrj [150] | Trajectory-level data for training universal force fields | Forces, atomic positions, energies | MACE-MP-0, NequIP, Allegro | Open-source |
| Pymatgen [174] | Tool for processing, analyzing, and visualizing material data | Crystal structures, phase diagrams, properties | Various | Open-source |
| $M^2$Hub [175] | Tool for data processing, model implementation, and training | Various | Various | Open-source |
| FireWorks [176] | Tool for end-to-end scientific workflows | Various | Job monitor, databases, high-performance computing systems | Open-source |
| Maggma [177] | Tool for data queries, transformations, and storage | Various | Databases | Open-source |
| ToolHub [36] | Unified interface with several tools to access online sources | Text | HoneyComb | Open-source |
| **Model development tools** | | | | |
| MatBench [178] | Benchmarking model performance across 13 materials tasks | Structure, composition, scalar properties | ALIGNN, MEGNet, Roost, CGCNN | Open-source |
| ALIGNN-FF (Toolkit) [179] | Training force fields from trajectories using message passing | Forces, energies, atomic trajectories | ALIGNN, ALIGNN-FF | Open-source |
| FORGE [43] | Pretraining and fine-tuning for foundation models | Graphs, atomic structures, text | ChemGPT, ALIGNN, MoLFormer variants | Open-source |
| OC20 (Tasks) [104] | Large-scale catalyst surface dataset and ML challenge suite | Adsorbates, surfaces, trajectories, energies | DimeNet++, SchNet, GemNet-OC | Open-source |
| Open MatSci ML Toolkit [42] | Standardized access to datasets, training and benchmarking | Structure, composition, text | ALIGNN, CGCNN, MEGNet | Open-source |
| A-Lab (Autonomous Lab) [26] | Robotic synthesis and closed-loop materials optimization | Protocols, text, synthesis metadata | Surrogate models, active learners | Closed-lab |
| GNoME Infrastructure (Discovery Engine) [14] | Active learning and scalable DFT validation for discovery | Graphs, formation energy, symmetry | GNN ensemble, GNoME | Internal |
| MatterSim Infrastructure (HPC Platform) [15] | Large-scale pretraining on 17M DFT-labeled structures | Atomic structures, forces, temperature | MatterSim | Internal |
| Ollama [180] | LLM deployment tool | Text | Various LLMs | Open-source |
| LangChain [181] | LLM-driven application development tool | Various | LLMs, databases, web scrapping | Open-source |

for exchanging material-related data, **OPTIMADE** [173] is an open-source tool for providing access to approximately 59M structures from 29 databases.

### 4.2.2  Model development tools

On the infrastructure and tooling front, several modular and open-source tools have emerged to support training, fine-tuning, and evaluation of both deep learning and foundation models. Focusing on the evaluation of machine learning and deep learning models trained for materials science, **MatBench** [178] leverages materials data from the aforementioned datasets, like MP, to formulate a comprehensive set of benchmarking tasks. **FORGE** [43] provides a flexible pretraining and benchmarking framework for large-scale scientific models, supporting graph-based architectures and self-supervised objectives across chemistry and materials datasets. Autonomous laboratories such as **A-Lab** [26] further illustrate the integration of foundation models with experimental loops, enabling robotic synthesis, online decision-making, and self-improving discovery. **Open MatSci ML Toolkit** [42] offers a curated interface to datasets, model architectures, and evaluation protocols, aimed at improving reproducibility and benchmarking across tasks. Toolkits like **ALIGNN-FF** [179] extend message-passing neural networks to force field learning, with a focus on universal atomistic simulation capabilities. **Geom3D** is a platform for geometric modeling on 3D structures and integrates MatBench and QMOF datasets and several geometry graphical neural networks. **JARVIS-Leaderboard** is an open-source and community-driven platform and allows users to set up benchmarks with various types of input data and different tasks. **GT4SD** [25] is a Python toolkit for developing and training generative models for scientific discovery. Training large-scale models also necessitates extensive infrastructures or resources. For example, **MatterSim** [15] is trained on over 17M DFT-labeled structures using industrial-scale infrastructures, while **GNoME** [14] leverages active learning in combination with distributed DFT validation. These efforts highlight the importance of scalable computational resources and data workflows in developing large foundation models.

The development of LLM and LLM-based agents requires the integration of LLM deployment, environment interaction, and external resource access. Several toolkits have been introduced to support these tasks. For instance, **Ollama** [180] provides a platform to deploy LLMs in local machines, covering a diverse set of open-source pretrained models. **LangChain** [181] is a framework for developing LLM-based applications, such as Retrieval Augmented Generation (RAG) or LLM agents. This framework provides the flexibility to connect different components and integrate third-party applications into the unified AI system. Similarly, **AutoGen** [186] and **CrewAI** [187] are popular frameworks for building role-playing LLM agents, facilitating the development of multi-agent systems. To support the deployment of an LLM-based agentic system which connects LLM to external data resources, **LlamaIndex** [188] is specifically introduced to provide the solution for this task.

Despite the scale of these efforts, access remains uneven. Most large-scale pretraining runs rely on proprietary or institutional resources, limiting reproducibility and broader community participation. Continued development of open infrastructures, shared checkpoints, and model zoos—alongside high-quality dataset curation—will be essential to democratize foundation model development and foster inclusive progress across materials science. Beyond static tooling, some models integrate tightly with experimental and high-performance computing pipelines.

## 5   Successes, Limitations, Challenges, and Future

With the recent rapid development of Artificial Intelligence, its application to scientific domains such as materials science is inevitable. AI offers predictive capabilities, insightful analysis, scalability, and efficiency that can significantly advance materials research. While the materials science community has been developing and utilizing datasets and programming toolkits for decades, foundation models are a more recent addition. Figure 3, which charts the release years of foundation models, datasets, and infrastructure tools, illustrates the accelerating convergence of AI and materials science. In this section, we first review the early successes of AI in the field, then examine current limitations and challenges, and finally discuss promising future research directions.
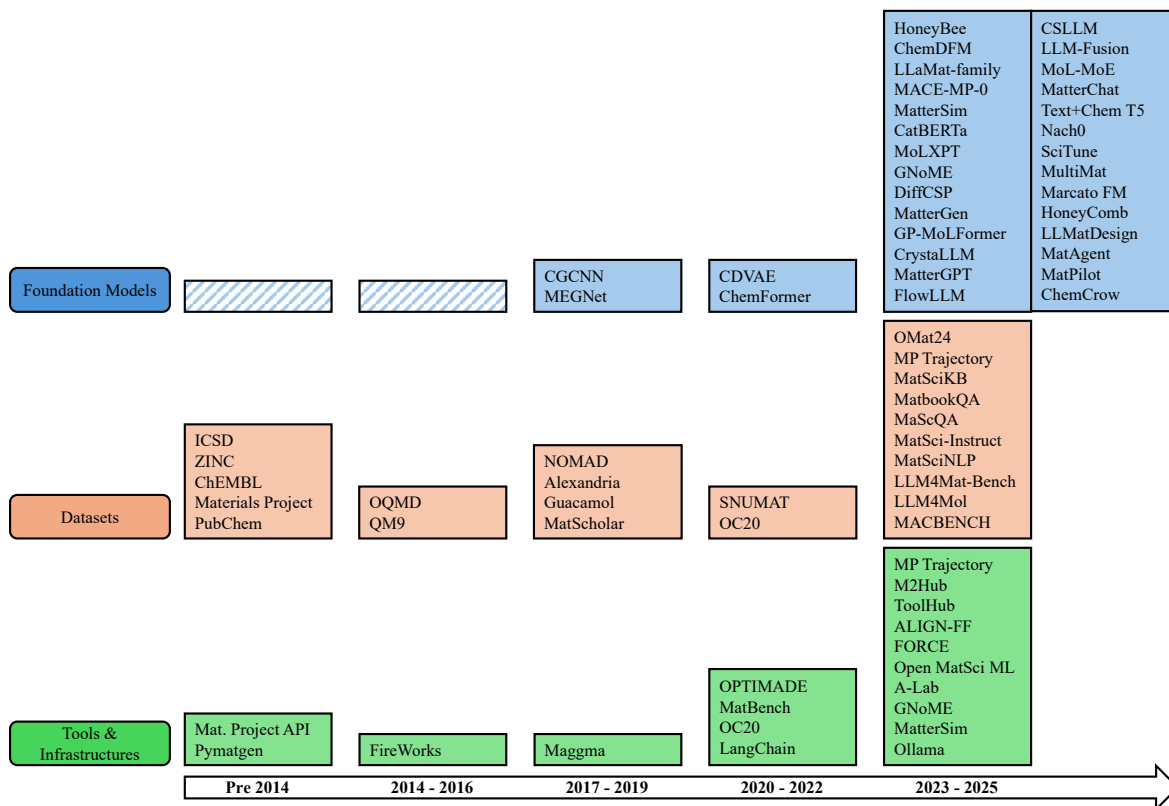


Figure 3: Development of AI in materials science over time: foundation models, datasets, and tools and infrastructure.

## 5.1  Early Successes

The growing maturity of foundation models in materials science is reflected in a series of high-impact applications that demonstrate their ability to augment or outperform traditional computational and experimental techniques. These successes span materials discovery, molecular dynamics, inverse design, autonomous synthesis, and literature-based information extraction, catalyzing a shift toward generalist, scalable, and efficient AI systems for scientific research [14, 15, 17, 20, 26].

One of the most significant achievements is the large-scale discovery of new materials. The GNoME (Graph Networks for Materials Exploration) system exemplifies this potential by using ensembles of graph neural networks trained on DFT formation energy data, symmetry-aware crystal generation, and iterative DFT validation in an active learning loop. GNoME identified over 2.2M new stable inorganic materials—an order-of-magnitude leap beyond prior efforts—and discovered over 45,000 new crystal prototypes occupying previously unexplored regions of chemical space [14]. Meanwhile, MatterSim, a universal machine-learned interatomic potential trained on over 17M DFT-labeled configurations, enables zero-shot molecular dynamics simulations with energy prediction errors below 50 meV/atom, even in high-temperature and non-equilibrium regimes. Other models such as MACE-MP-0 have demonstrated near-DFT accuracy using message-passing neural networks applied to periodic systems.

In generative modeling and inverse design, diffusion- and language-based foundation models are enabling property-guided exploration of materials space. MatterGen combines unsupervised pretraining with adapter-based fine-tuning to generate chemically valid, structurally stable, and property-optimized inorganic crystals. It outperforms prior models in diversity, novelty, and closeness to DFT-derived properties. Language-based models like GP-MoLFormer extend molecular design to functional materials and drug-like compounds through property-conditioned generation.

On the text mining and scientific reasoning front, models like nach0 and ChemDFM demonstrate cross-modal capabilities for molecule generation, retrosynthesis, property prediction, and question answering. Trained on structured and unstructured chemistry corpora, these models extract synthesis routes, interpret SMILES strings, and enable high-throughput knowledge extraction from literature.

Although LLM-based agents are still in their early stages and continue to evolve toward more sophisticated and practical applications, recent developments—such as HoneyComb, LLMatDesign, MatAgent, and MatPilot—demonstrate promising progress in building agentic systems for materials science. Beyond prediction and generation, several of these works incorporate foundation models into agentic workflows that support planning, reasoning, and tool integration. These systems offer important conceptual frameworks that aid researchers and practitioners in understanding, analyzing, and automating complex tasks more efficiently. For example, MatAgent, LLMatDesign, and MatPilot utilize LLMs to coordinate multi-step processes such as candidate screening, simulation planning, and synthesis route generation, all with minimal human intervention.

Foundation models have also been deployed in autonomous systems: the A-Lab integrates language planning, thermodynamic reasoning, robotics, and active learning into a closed-loop experimental platform. Operating over 17 days, A-Lab synthesized 41 of 58 targeted materials with a 71% success rate and minimal human intervention.

## 5.2  Limitations and Challenges

Despite recent promising developments, foundation models in materials science remain an emerging technology with substantial limitations. One of the most critical challenges is the accurate modeling of long-range interactions such as electrostatics, dispersion, and magnetism. Local message-passing architectures like those in MatterSim and MACE-MP-0 effectively capture short-range bonding but fall short for systems dominated by nonlocal physics [16, 15]. Although empirical corrections can partially mitigate this issue, comprehensive and efficient solutions remain elusive.

Generalizability presents a significant challenge in the training of foundation models for materials science. A core objective in this field is the discovery of novel materials with unprecedented properties tailored for specific applications. Foundation models trained on existing material distributions should demonstrate the capability to be adaptive to out-of-distribution materials or unseen domains. Moreover, materials behave differently under different conditions, such as super high/low temperatures or pressure. Vision-based foundation models, used to detect defects from images or videos, may also struggle as the visual appearance of materials evolves over time. In addition, many underrepresented material classes—such as polymers, disordered solids, and biomaterials—remain challenging due to sparse, noisy, or long-tail data distributions, which limit the generalizability of current FMs. These challenges underscore the need for more in-depth research into the design and training of robust foundation models in materials science.

Integration with experimental workflows remains limited. Most foundation models are not trained with synthesizability, phase stability, or safety in mind. Although MatterGen can generate candidate materials optimized for electronic or mechanical properties, it cannot guarantee experimental feasibility [17]. Few models jointly optimize for both

computational performance and laboratory viability. Similarly, multimodal integration remains underdeveloped. While models like nach0 and ATLANTIC explore text-structure-property fusion, most foundation models operate in unimodal spaces [20, 23], lacking the ability to simultaneously reason over images, spectra, simulation outputs, and experimental metadata.

Interpretability is another pressing concern. As these models increase in size and complexity, understanding their internal representations and assessing their trustworthiness becomes increasingly difficult. Generative models like nach0 and MatterGen can produce physically implausible results despite being syntactically valid [20, 17]. Black-box behavior poses risks in high-stakes applications like battery chemistry or catalysis. Furthermore, foundation models are often evaluated on in-distribution test sets, whereas materials discovery requires extrapolation to new chemistries, properties, or structural motifs. While models like GNoME have demonstrated generalization across broad chemical spaces, the reliability of predictions in entirely novel regimes is unclear [14].

Data bias is another significant limitation. Current training corpora overrepresent stable, inorganic, equilibrium-phase systems—especially oxides—while underrepresenting high-entropy alloys, amorphous phases, soft materials, and molecular crystals. While active learning and off-equilibrium sampling strategies aim to address this, coverage remains incomplete. This imbalance reduces the generality of models and limits their utility in less-studied or industrially relevant materials domains.

Several challenges arise in the design and deployment of LLMs and LLM agents, necessitating further research. LLMs rely heavily on the availability of data to train the language models. The data collection and preprocessing stages are critical for ensuring the quality of the generated content, particularly when natural and human-like responses are expected. These stages often require significant resources and costs to curate high-quality training data.

Another key challenge is LLM hallucination, especially when models are applied to unfamiliar or out-of-distribution scenarios. Such issues may stem from corrupted training data, limitations in model architecture, or ambiguous user instructions. This might arise due to the corrupted training data, limited model design, or ambiguous user instructions. Finally, biosafety and chemical safety are highly important while working with materials science experiments. These models may generate syntactically valid but physically implausible results, or propose unsafe synthesis conditions without proper grounding. Experimental plans generated by LLM agents might offer excellent opportunities to analyze and discover new findings but also contain potential risks if not thoroughly reviewed by human experts. These considerations highlight the urgent need for comprehensive evaluation benchmarks and training frameworks (e.g. supervised fine-tuning, reinforcement learning, human-in-the-loop) tailored to specific tasks in materials science, to ensure the safe and effective deployment of LLM and LLM agents.

Supporting the development of autonomous, LLM-based agentic systems also requires the creation of open-source, highly scalable, and integrable materials science tools and infrastructure. High-quality, standardized tools can significantly reduce the time and resources required for implementation, thereby accelerating the transition of agentic systems from research prototypes to practical, real-world applications accessible to researchers and end users.

The substantial compute requirements for training and deploying foundation models raise significant concerns about accessibility. For example, models such as MatterSim and GNoME demand tens of thousands of GPU hours, while large language models like GPT, Gemini, and LLaMA require millions of GPU hours, often relying on proprietary computing resources—placing such efforts beyond the reach of most academic researchers. Open-source tools and frameworks like FORGE, Open MatSci ML, and ALIGNN-FF represent promising steps toward democratizing access and enabling smaller-scale experimentation. Nevertheless, despite these advances, fostering broader participation and ensuring sustainable development across the research community remain ongoing challenges that must be addressed to fully realize the potential of foundation models in materials science.

### 5.3 Future Direction

To address these limitations and chart the path forward, several directions merit emphasis. Future foundation models must incorporate stronger physical laws to align with quantum mechanics, thermodynamics, and symmetry principles. Architectures with built-in equivariance or physics-informed constraints may improve both accuracy and interpretability. Cross-modal models that integrate crystal structures, synthesis routes, spectra, images, and text hold the key to more holistic materials reasoning. Expanding training datasets to include failed experiments, synthesis outcomes, and non-traditional materials classes will improve generality and real-world relevance. Another direction is the principled development of foundation models for materials science, which are typically data hungry, that can assimilate both low-fidelity data (cheap to collect) and high-fidelity data (expensive to collect) to circumvent the high cost of building a large high-fidelity training dataset and mitigate data sparsity concerns. Progress will also depend on aligning text, structure, spectra, and image-based data through high-quality, co-registered datasets and multimodal pretraining strategies. Unified tokenization schemes and modality adapters may be necessary to bridge diverse formats.

Active and continual learning will be essential for deploying foundation models in dynamic, exploratory settings. Integration with robotic labs, uncertainty-aware retraining, and human-in-the-loop frameworks will allow models to learn from failure, adapt to new tasks, and serve as collaborators in discovery pipelines. Efficient training paradigms—modular, sparse, or low-rank—will be necessary to reduce the environmental and financial cost of model development. The growing potential for combining large language models like GPT-4 with specialized materials tools signals a future where AI-driven autonomous agents could revolutionize materials science research, making these models indispensable to materials discovery. Alongside technical progress, the field must prioritize open benchmarks, reproducible evaluation, interpretability tooling, and shared infrastructures to ensure equitable progress.

Data governance and data curation are essential steps in enhancing the quality of training data. Human feedback data is also a crucial component in achieving more natural conversations between users and models. The principles of trustworthy AI can be employed to develop models that prioritize privacy, safety, and reliability. Techniques such as adversarial learning and machine unlearning contribute to these goals by improving model robustness and enabling the removal of sensitive or erroneous data. Finally, integrating human-in-the-loop approaches offers a promising strategy for incorporating human oversight during task execution, thereby improving both accuracy and accountability.

Trustworthy AI principles such as uncertainty quantification, adversarial testing, machine unlearning, and human feedback must guide FM deployment in scientific domains. Ensuring safe, robust, and interpretable behavior will be essential for adoption in high-stakes applications. Ultimately, the success of foundation models in materials science hinges on interdisciplinary collaboration, standardized data practices, and deep integration of domain knowledge and experimental feedback. With these components in place, foundation models have the potential to become not just predictive engines, but collaborative agents for reasoning, creativity, and exploration across materials science.

# 6    Conclusion

The emergence of foundation models represents a profound shift in materials science, echoing the transformations currently underway in fields like natural language processing, computer vision, and molecular biology. These large, pretrained, and general-purpose models have demonstrated remarkable capabilities in property prediction, atomistic simulation, materials generation, and language-based scientific reasoning. What distinguishes foundation models from previous generations of machine learning tools is not only their scale, but also their transferability, emergent behavior, and versatility across tasks and data modalities.

In this survey, we propose a clear taxonomy of materials foundation models, organizing them according to tasks—including data extraction, atomistic simulation, property prediction, materials design and discovery, process optimization, and multiscale modeling—and by model type, ranging from specialized unimodal foundation models to ambitious multimodal foundation models, as well as emerging autonomous and flexible LLM agents. Within these categories, unimodal foundation models such as GNoME, MatterSim, MACE-MP-0, and MatterGen collectively demonstrate the feasibility of AI-driven materials discovery at unprecedented scales and speeds. Multimodal foundation models like MatterChat, Text+Chem T5, and nach0 exemplify efforts to integrate multimodal materials information, thereby enhancing the quality of outcomes. The recent emergence of LLM-powered agents (e.g., HoneyComb, MatA-gent, LLMatDesign) highlights a shift toward systems that not only analyze data but also reason, plan, and interact across formats such as text, structure, images, and protocols.

While these successes are undeniably impressive, our review also surfaces significant limitations and open challenges. Issues of long-range interaction modeling, data imbalance, interpretability, and generalization under extrapolation remain unresolved. Similarly, the computational and infrastructural demands of training large foundation models limit their accessibility and reproducibility in the broader materials research community. Integrating these models into experimental workflows and achieving multimodal reasoning capabilities remain aspirational goals for the near term.

Looking ahead, the field is clearly moving toward scalable, multimodal, and human-AI collaborative discovery systems. Foundation models that can reason over atomic structures, properties, synthesis procedures, and scientific text within a unified framework have the potential to dramatically accelerate materials research, while also opening new frontiers in autonomous laboratories, hypothesis-driven design, and interactive scientific exploration. Realizing this vision will require not only continued advances in model architectures, datasets, and computational infrastructures, but also a cultural shift toward openness, interdisciplinarity, and rigorous, transparent evaluation. To realize their full potential, foundation models must be not only accurate and scalable, but also trustworthy, interpretable, and safe. Community efforts toward reproducibility, human-in-the-loop systems, governance, and equitable access will be critical in guiding their responsible deployment.

Ultimately, foundation models offer a powerful new paradigm for materials science; one that complements and enhances, rather than replaces, human expertise. The challenge for the research community now is to harness these tools

responsibly, inclusively, and creatively, ensuring that the benefits of AI-driven materials discovery are both scientifically rigorous and broadly accessible. Foundation models are poised to become collaborative scientific agents, transforming how we explore, reason about, and design the materials of the future.

## Acknowledgement

## References

[1] Edward O Pyzer-Knapp, Matteo Manica, Peter Staar, Lucas Morin, Patrick Ruch, Teodoro Laino, John R Smith, and Alessandro Curioni. Foundation models for materials discovery–current state and future directions. *npj Computational Materials*, 11(1):61, 2025.

[2] Zhenzhong Wang, Haowei Hua, Wanyu Lin, Ming Yang, and Kay Chen Tan. Crystalline material discovery in the era of artificial intelligence. *arXiv preprint arXiv:2408.08044*, 2024.

[3] Xiao-Qi Han, Xin-De Wang, Meng-Yuan Xu, Zhen Feng, Bo-Wen Yao, Peng-Jie Guo, Ze-Feng Gao, and Zhong-Yi Lu. Ai-driven inverse design of materials: Past, present and future. *Chinese Physics Letters*, 2024.

[4] Mayk Caldas Ramos, Christopher J Collison, and Andrew D White. A review of large language models and autonomous agents in chemistry. *Chemical Science*, 2025.

[5] Piyush Karande, Brian Gallagher, and Thomas Yong-Jin Han. A strategic approach to machine learning for material science: how to tackle real-world challenges and avoid pitfalls. *Chemistry of Materials*, 34(17): 7650–7665, 2022.

[6] Md Hosne Mobarak, Mariam Akter Mimona, Md Aminul Islam, Nayem Hossain, Fatema Tuz Zohura, Ibnul Imtiaz, and Md Israfil Hossain Rimon. Scope of machine learning in materials research—a review. *Applied Surface Science Advances*, 18:100523, 2023.

[7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 conference of the North American chapter of the association for computational linguistics: human language technologies, volume 1 (long and short papers)*, pages 4171–4186, 2019.

[8] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.

[9] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.

[10] Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altenschmidt, Sam Altman, Shyamal Anadkat, et al. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*, 2023.

[11] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113, 2023.

[12] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, pages 8748–8763. PmLR, 2021.

[13] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 9650–9660, 2021.

[14] Amil Merchant, Simon Batzner, Samuel S Schoenholz, Muratahan Aykol, Gowoon Cheon, and Ekin Dogus Cubuk. Scaling deep learning for materials discovery. *Nature*, 624(7990):80–85, 2023.

[15] Han Yang, Chenxi Hu, Yichi Zhou, Xixian Liu, Yu Shi, Jielan Li, Guanzhi Li, Zekun Chen, Shuizhou Chen, Claudio Zeni, et al. Mattersim: A deep learning atomistic model across elements, temperatures and pressures. *arXiv preprint arXiv:2405.04967*, 2024.

[16] Ilyes Batatia, Philipp Benner, Yuan Chiang, Alin M Elena, Dávid P Kovács, Janosh Riebesell, Xavier R Advincula, Mark Asta, Matthew Avaylon, William J Baldwin, et al. A foundation model for atomistic materials chemistry. *arXiv preprint arXiv:2401.00096*, 2023.

[17] Claudio Zeni, Robert Pinsler, Daniel Zügner, Andrew Fowler, Matthew Horton, Xiang Fu, Sasha Shysheya, Jonathan Crabbé, Lixin Sun, Jake Smith, et al. Mattergen: a generative model for inorganic materials design. *arXiv preprint arXiv:2312.03687*, 2023.

[18] Rui Jiao, Wenbing Huang, Yu Liu, Deli Zhao, and Yang Liu. Space group constrained crystal generation. *arXiv preprint arXiv:2402.03992*, 2024.

[19] Zhendong Cao, Xiaoshan Luo, Jian Lv, and Lei Wang. Space group informed transformer for crystalline materials generation. *arXiv preprint arXiv:2403.15734*, 2024.

[20] Micha Livne, Zulfat Miftahutdinov, Elena Tutubalina, Maksim Kuznetsov, Daniil Polykovskiy, Annika Brundyn, Aastha Jhunjhunwala, Anthony Costa, Alex Aliper, Alán Aspuru-Guzik, et al. nach0: multimodal natural and chemical languages foundation model. *Chemical Science*, 15(22):8380–8389, 2024.

[21] Viggo Moro, Charlotte Loh, Rumen Dangovski, Ali Ghorashi, Andrew Ma, Zhuo Chen, Samuel Kim, Peter Y Lu, Thomas Christensen, and Marin Soljačić. Multimodal foundation models for material property prediction and discovery. *Newton*, 2025.

[22] Yingheng Tang, Wenbin Xu, Jie Cao, Jianzhu Ma, Weilu Gao, Steve Farrell, Benjamin Erichson, Michael W Mahoney, Andy Nonaka, and Zhi Yao. Matterchat: A multi-modal llm for material science. *arXiv preprint arXiv:2502.13107*, 2025.

[23] Sai Munikoti, Anurag Acharya, Sridevi Wagle, and Sameera Horawalavithana. Atlantic: Structure-aware retrieval-augmented language model for interdisciplinary science. *arXiv preprint arXiv:2311.12289*, 2023.

[24] Luis M Antunes, Keith T Butler, and Ricardo Grau-Crespo. Crystal structure generation with autoregressive large language modeling. *Nature Communications*, 15(1):1–16, 2024.

[25] Matteo Manica, Jannis Born, Joris Cadow, Dimitrios Christofidellis, Ashish Dave, Dean Clarke, Yves Gaetan Nana Teukam, Giorgio Giannone, Samuel C Hoffman, Matthew Buchan, et al. Accelerating material design with the generative toolkit for scientific discovery. *npj Computational Materials*, 9(1):69, 2023.

[26] Nathan J Szymanski, Bernardus Rendy, Yuxing Fei, Rishi E Kumar, Tanjin He, David Milsted, Matthew J McDermott, Max Gallant, Ekin Dogus Cubuk, Amil Merchant, et al. An autonomous laboratory for the accelerated synthesis of novel materials. *Nature*, 624(7990):86–91, 2023.

[27] Agnese Marcato, Javier E Santos, Aleksandra Pachalieva, Kai Gao, Ryley Hill, Esteban Rougier, Qinjun Kang, Jeffrey Hyman, Abigail Hunter, Janel Chua, et al. Developing a foundation model for predicting material failure. *arXiv preprint arXiv:2411.08354*, 2024.

[28] Lei Ren, Haiteng Wang, Yuqing Wang, Keke Huang, Lihui Wang, and Bohu Li. Foundation models for the process industry: Challenges and opportunities. *Engineering*, 2025.

[29] Junyoung Choi, Gunwook Nam, Jaesik Choi, and Yousung Jung. A perspective on foundation models in chemistry. *JACS Au*, 5(4):1499–1518, 2025.

[30] Kin Long Kelvin Lee, Carmelo Gonzales, Matthew Spellings, Mikhail Galkin, Santiago Miret, and Nalini Kumar. Towards foundation models for materials science: The open matsci ml toolkit. In *Proceedings of the SC'23 Workshops of the International Conference on High Performance Computing, Network, Storage, and Analysis*, pages 51–59, 2023.

[31] Kamal Choudhary. Atomgpt: Atomistic generative pretrained transformer for forward and inverse materials design. *The Journal of Physical Chemistry Letters*, 15(27):6909–6917, 2024.

[32] Eduardo Soares, Indra Priyadarsini, Emilio Vital Brazil, Victor Yukio Shirasuna, and Seiji Takeda. Multi-view mixture-of-experts for predicting molecular properties using smiles, selfies, and graph-based representations. In *Neurips 2024 Workshop Foundation Models for Science: Progress, Opportunities, and Challenges*, 2024.

[33] Zihan Zhao, Da Ma, Lu Chen, Liangtai Sun, Zihao Li, Yi Xia, Bo Chen, Hongshen Xu, Zichen Zhu, Su Zhu, et al. Chemdfm: A large language foundation model for chemistry. *arXiv preprint arXiv:2401.14818*, 2024.

[34] Vaibhav Mishra, Somaditya Singh, Dhruv Ahlawat, Mohd Zaki, Vaibhav Bihani, Hargun Singh Grover, Biswajit Mishra, Santiago Miret, NM Krishnan, et al. Foundational large language models for materials research. *arXiv preprint arXiv:2412.09560*, 2024.

[35] Sameera Horawalavithana, Sai Munikoti, Ian Stewart, and Henry Kvinge. Scitune: Aligning large language models with scientific multimodal instructions. *arXiv preprint arXiv:2307.01139*, 2023.

[36] Huan Zhang, Yu Song, Ziyu Hou, Santiago Miret, and Bang Liu. Honeycomb: A flexible llm-based agent system for materials science. *arXiv preprint arXiv:2409.00135*, 2024.

[37] Shuyi Jia, Chao Zhang, and Victor Fung. Llmatdesign: Autonomous materials discovery with large language models. *arXiv preprint arXiv:2406.13163*, 2024.

[38] Yeonghun Kang and Jihan Kim. Chatmof: an artificial intelligence system for predicting and generating metal-organic frameworks using large language models. *Nature communications*, 15(1):4705, 2024.

[39] Adib Bazgir, Yuwen Zhang, et al. Matagent: A human-in-the-loop multi-agent llm framework for accelerating the material science discovery cycle. In *AI for Accelerated Materials Design-ICLR 2025*.

[40] Ziqi Ni, Yahao Li, Kaijia Hu, Kunyuan Han, Ming Xu, Xingyu Chen, Fengqi Liu, Yicong Ye, and Shuxin Bai. Matpilot: an llm-enabled ai materials scientist under the framework of human-machine collaboration. *arXiv preprint arXiv:2411.08063*, 2024.

[41] Andres M. Bran, Sam Cox, Oliver Schilter, Carlo Baldassari, Andrew D White, and Philippe Schwaller. Augmenting large language models with chemistry tools. *Nature Machine Intelligence*, 6(5):525–535, 2024.

[42] Santiago Miret, Kin Long Kelvin Lee, Carmelo Gonzales, Marcel Nassar, and Matthew Spellings. The open matsci ml toolkit: A flexible framework for machine learning in materials science. *arXiv preprint arXiv:2210.17484*, 2022.

[43] Junqi Yin, Sajal Dash, Feiyi Wang, and Mallikarjun Shankar. Forge: Pre-training open foundation models for science. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–13, 2023.

[44] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.

[45] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. Improving language understanding by generative pre-training. 2018.

[46] Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-Baptiste Alayrac, Jiahui Yu, Radu Soricut, Johan Schalkwyk, Andrew M Dai, Anja Hauth, Katie Millican, et al. Gemini: a family of highly capable multimodal models. *arXiv preprint arXiv:2312.11805*, 2023.

[47] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.

[48] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*, 2019.

[49] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. In *International Conference on Learning Representations (ICLR)*, 2023.

[50] Chan Hee Song, Jiaman Wu, Clayton Washington, Brian M Sadler, Wei-Lun Chao, and Yu Su. Llm-planner: Few-shot grounded planning for embodied agents with large language models. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 2998–3009, 2023.

[51] Shunyu Yao, Dian Yu, Jeffrey Zhao, Izhak Shafran, Tom Griffiths, Yuan Cao, and Karthik Narasimhan. Tree of thoughts: Deliberate problem solving with large language models. *Advances in neural information processing systems*, 36:11809–11822, 2023.

[52] Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents. In *International conference on machine learning*, pages 9118–9147. PMLR, 2022.

[53] Michal Kosinski. Theory of mind may have spontaneously emerged in large language models. *arXiv preprint arXiv:2302.02083*, 4:169, 2023.

[54] Murray Shanahan, Kyle McDonell, and Laria Reynolds. Role play with large language models. *Nature*, 623 (7987):493–498, 2023.

[55] Lei Wang, Chen Ma, Xueyang Feng, Zeyu Zhang, Hao Yang, Jingsen Zhang, Zhiyuan Chen, Jiakai Tang, Xu Chen, Yankai Lin, et al. A survey on large language model based autonomous agents. *Frontiers of Computer Science*, 18(6):186345, 2024.

[56] Yu Song, Santiago Miret, Huan Zhang, and Bang Liu. Honeybee: Progressive instruction finetuning of large language models for materials science. *arXiv preprint arXiv:2310.08511*, 2023.

[57] Tian Xie and Jeffrey C Grossman. Crystal graph convolutional neural networks for an accurate and interpretable prediction of material properties. *Physical review letters*, 120(14):145301, 2018.

[58] Chi Chen, Weike Ye, Yunxing Zuo, Chen Zheng, and Shyue Ping Ong. Graph networks as a universal machine learning framework for molecules and crystals. *Chemistry of Materials*, 31(9):3564–3572, 2019.

[59] Janghoon Ock, Chakradhar Guntuboina, and Amir Barati Farimani. Catalyst energy prediction with catberta: Unveiling feature exploration strategies through large language models. *ACS Catalysis*, 13(24):16032–16044, 2023.

[60] Zequn Liu, Wei Zhang, Yingce Xia, Lijun Wu, Shufang Xie, Tao Qin, Ming Zhang, and Tie-Yan Liu. Molxpt: Wrapping molecules with text for generative pre-training. *arXiv preprint arXiv:2305.10688*, 2023.

[61] Tian Xie, Xiang Fu, Octavian-Eugen Ganea, Regina Barzilay, and Tommi Jaakkola. Crystal diffusion variational autoencoder for periodic material generation. *arXiv preprint arXiv:2110.06197*, 2021.

[62] Rui Jiao, Wenbing Huang, Peijia Lin, Jiaqi Han, Pin Chen, Yutong Lu, and Yang Liu. Crystal structure prediction by joint equivariant diffusion. *Advances in Neural Information Processing Systems*, 36:17464–17497, 2023.

[63] Jerret Ross, Brian Belgodere, Samuel C Hoffman, Vijil Chenthamarakshan, Jiri Navratil, Youssef Mroueh, and Payel Das. Gp-molformer: A foundation model for molecular generation. *arXiv preprint arXiv:2405.04912*, 2024.

[64] Yan Chen, Xueru Wang, Xiaobin Deng, Yilun Liu, Xi Chen, Yunwei Zhang, Lei Wang, and Hang Xiao. Mattergpt: A generative transformer for multi-property inverse design of solid-state materials. *arXiv preprint arXiv:2408.07608*, 2024.

[65] Ross Irwin, Spyridon Dimitriadis, Jiazhen He, and Esben Jannik Bjerrum. Chemformer: a pre-trained transformer for computational chemistry. *Machine Learning: Science and Technology*, 3(1):015022, 2022.

[66] Anuroop Sriram, Benjamin Miller, Ricky TQ Chen, and Brandon Wood. Flowllm: Flow matching for material generation with large language models as base distributions. *Advances in Neural Information Processing Systems*, 37:46025–46046, 2024.

[67] Zhilong Song, Shuaihua Lu, Minggang Ju, Qionghua Zhou, and Jinlan Wang. Is large language model all you need to predict the synthesizability and precursors of crystal structures? *arXiv preprint arXiv:2407.07016*, 2024.

[68] Onur Boyar, Indra Priyadarsini, Seiji Takeda, and Lisa Hamada. Llm-fusion: A novel multimodal fusion model for accelerated material discovery. *arXiv preprint arXiv:2503.01022*, 2025.

[69] Dimitrios Christofidellis, Giorgio Giannone, Jannis Born, Ole Winther, Teodoro Laino, and Matteo Manica. Unifying molecular and textual representations via multi-task language modelling. In *International Conference on Machine Learning*, pages 6140–6157. PMLR, 2023.

[70] Amalie Trewartha, Nicholas Walker, Haoyan Huo, Sanghoon Lee, Kevin Cruse, John Dagdelen, Alexander Dunn, Kristin A Persson, Gerbrand Ceder, and Anubhav Jain. Quantifying the advantage of domain-specific pre-training on named entity recognition tasks in materials science. *Patterns*, 3(4), 2022.

[71] Tanishq Gupta, Mohd Zaki, NM Anoop Krishnan, and Mausam. Matscibert: A materials domain language model for text mining and information extraction. *npj Computational Materials*, 8(1):102, 2022.

[72] Pranav Shetty, Arunkumar Chitteth Rajan, Chris Kuenneth, Sonakshi Gupta, Lakshmi Prerana Panchumarti, Lauren Holm, Chao Zhang, and Rampi Ramprasad. A general-purpose material property data extraction pipeline from large polymer corpora using natural language processing. *npj Computational Materials*, 9(1):52, 2023.

[73] Yujie Qian, Jiang Guo, Zhengkai Tu, Zhening Li, Connor W Coley, and Regina Barzilay. Molscribe: robust molecular structure recognition with image-to-graph generation. *Journal of Chemical Information and Modeling*, 63(7):1925–1934, 2023.

[74] Fangyu Liu, Julian Martin Eisenschlos, Francesco Piccinno, Syrine Krichene, Chenxi Pang, Kenton Lee, Mandar Joshi, Wenhu Chen, Nigel Collier, and Yasemin Altun. Deplot: One-shot visual language reasoning by plot-to-table translation. *arXiv preprint arXiv:2212.10505*, 2022.

[75] Mohd Zaki, NM Anoop Krishnan, et al. Mascqa: investigating materials science knowledge of large language models. *Digital Discovery*, 3(2):313–327, 2024.

[76] Joren Van Herck, María Victoria Gil, Kevin Maik Jablonka, Alex Abrudan, Andy S Anker, Mehrdad Asgari, Ben Blaiszik, Antonio Buffo, Leander Choudhury, Clemence Corminboeuf, et al. Assessment of fine-tuned large language models for real-world chemistry and material science applications. *Chemical science*, 16(2):670–684, 2025.

[77] Ben Wang and Aran Komatsuzaki. Gpt-j-6b: A 6 billion parameter autoregressive language model, 2021.

[78] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.

[79] Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. Mistral 7b, 2023. URL `https://arxiv.org/abs/2310.06825`.

[80] David Weininger. Smiles, a chemical language and information system. 1. introduction to methodology and encoding rules. *Journal of chemical information and computer sciences*, 28(1):31–36, 1988.

[81] Hengxing Cai, Xiaochen Cai, Shuwen Yang, Jiankun Wang, Lin Yao, Zhifeng Gao, Junhan Chang, Sihang Li, Mingjun Xu, Changxin Wang, et al. Uni-smart: Universal science multimodal analysis and research transformer. *arXiv preprint arXiv:2403.10301*, 2024.

[82] Anthropic. The claude 3 model family: Opus, sonnet, haiku. 2025. URL `https://api.semanticscholar.org/CorpusID:268232499`.

[83] Chengxuan Ying, Tianle Cai, Shengjie Luo, Shuxin Zheng, Guolin Ke, Di He, Yanming Shen, and Tie-Yan Liu. Do transformers really perform badly for graph representation? *Advances in neural information processing systems*, 34:28877–28888, 2021.

[84] Ilyes Batatia, David P Kovacs, Gregor Simm, Christoph Ortner, and Gábor Csányi. Mace: Higher order equivariant message passing neural networks for fast and accurate force fields. *Advances in neural information processing systems*, 35:11423–11436, 2022.

[85] Justin S Smith, Olexandr Isayev, and Adrian E Roitberg. Ani-1: an extensible neural network potential with dft accuracy at force field computational cost. *Chemical science*, 8(4):3192–3203, 2017.

[86] Jörg Behler and Michele Parrinello. Generalized neural-network representation of high-dimensional potential-energy surfaces. *Physical review letters*, 98(14):146401, 2007.

[87] Roman Zubatyuk, Justin S Smith, Jerzy Leszczynski, and Olexandr Isayev. Accurate and transferable multitask prediction of chemical properties with an atoms-in-molecules neural network. *Science advances*, 5(8):eaav6490, 2019.

[88] Dylan M Anstine, Roman Zubatyuk, and Olexandr Isayev. Aimnet2: a neural network potential to meet your neutral, charged, organic, and elemental-organic needs. *Chemical Science*, 2025.

[89] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. The graph neural network model. *IEEE transactions on neural networks*, 20(1):61–80, 2008.

[90] Mario Krenn, Florian Häse, AkshatKumar Nigam, Pascal Friederich, and Alan Aspuru-Guzik. Self-referencing embedded strings (selfies): A 100% robust molecular string representation. *Machine Learning: Science and Technology*, 1(4):045024, 2020.

[91] Teague Sterling and John J Irwin. Zinc 15–ligand discovery for everyone. *Journal of chemical information and modeling*, 55(11):2324–2337, 2015.

[92] John J Irwin, Khanh G Tang, Jennifer Young, Chinzorig Dandarchuluun, Benjamin R Wong, Munkhzul Khurel-baatar, Yurii S Moroz, John Mayfield, and Roger A Sayle. Zinc20—a free ultralarge-scale chemical database for ligand discovery. *Journal of chemical information and modeling*, 60(12):6065–6073, 2020.

[93] Barbara Zdrazil, Eloy Felix, Fiona Hunter, Emma J Manners, James Blackshaw, Sybilla Corbett, Marleen de Veij, Harris Ioannidis, David Mendez Lopez, Juan F Mosquera, et al. The chembl database in 2023: a drug discovery platform spanning multiple bioactivity data types and time periods. *Nucleic acids research*, 52(D1): D1180–D1192, 2024.

[94] Mark Davies, Michał Nowotka, George Papadatos, Nathan Dedman, Anna Gaulton, Francis Atkinson, Louisa Bellis, and John P Overington. Chembl web services: streamlining access to drug discovery data and utilities. *Nucleic acids research*, 43(W1):W612–W620, 2015.

[95] Kristof T Schütt, Huziel E Sauceda, P-J Kindermans, Alexandre Tkatchenko, and K-R Müller. Schnet–a deep learning architecture for molecules and materials. *The Journal of Chemical Physics*, 148(24), 2018.

[96] Anubhav Jain, Shyue Ping Ong, Geoffroy Hautier, Wei Chen, William Davidson Richards, Stephen Dacek, Shreyas Cholia, Dan Gunter, David Skinner, Gerbrand Ceder, et al. Commentary: The materials project: A materials genome approach to accelerating materials innovation. *APL materials*, 1(1), 2013.

[97] Lars Ruddigkeit, Ruud Van Deursen, Lorenz C Blum, and Jean-Louis Reymond. Enumeration of 166 billion organic small molecules in the chemical universe database gdb-17. *Journal of chemical information and modeling*, 52(11):2864–2875, 2012.

[98] Raghunathan Ramakrishnan, Pavlo O Dral, Matthias Rupp, and O Anatole Von Lilienfeld. Quantum chemistry structures and properties of 134 kilo molecules. *Scientific data*, 1(1):1–7, 2014.

[99] Oliver Wieder, Stefan Kohlbacher, Mélaine Kuenemann, Arthur Garon, Pierre Ducrot, Thomas Seidel, and Thierry Langer. A compact review of molecular property prediction with graph neural networks. *Drug Discovery Today: Technologies*, 37:1–12, 2020.

[100] Seyone Chithrananda, Gabriel Grand, and Bharath Ramsundar. Chemberta: large-scale self-supervised pretraining for molecular property prediction. *arXiv preprint arXiv:2010.09885*, 2020.

[101] Juncai Li and Xiaofei Jiang. Mol-bert: An effective molecular representation with bert for molecular property prediction. *Wireless Communications and Mobile Computing*, 2021(1):7181815, 2021.

[102] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.

[103] Sunghwan Kim, Jie Chen, Tiejun Cheng, Asta Gindulyte, Jia He, Siqian He, Qingliang Li, Benjamin A Shoemaker, Paul A Thiessen, Bo Yu, et al. Pubchem 2025 update. *Nucleic Acids Research*, 53(D1):D1516–D1525, 2025.

[104] Lowik Chanussot, Abhishek Das, Siddharth Goyal, Thibaut Lavril, Muhammed Shuaibi, Morgane Riviere, Kevin Tran, Javier Heras-Domingo, Caleb Ho, Weihua Hu, et al. Open catalyst 2020 (oc20) dataset and community challenges. *Acs Catalysis*, 11(10):6059–6072, 2021.

[105] Atakan Yüksel, Erva Ulusoy, Atabey Ünlü, and Tunca Doğan. Selformer: molecular representation learning via selfies language models. *Machine Learning: Science and Technology*, 4(2):025035, 2023.

[106] Keqiang Yan, Yi Liu, Yuchao Lin, and Shuiwang Ji. Periodic graph transformers for crystal material property prediction. *Advances in Neural Information Processing Systems*, 35:15066–15080, 2022.

[107] Keqiang Yan, Cong Fu, Xiaofeng Qian, Xiaoning Qian, and Shuiwang Ji. Complete and efficient graph transformers for crystal material property prediction. *arXiv preprint arXiv:2403.11857*, 2024.

[108] Sanjar Adilov. Generative pre-training from molecules. 2021.

[109] National Library of Medicine (US). PubMed Database. `https://pubmed.ncbi.nlm.nih.gov/`, 2024. Accessed: 2025-06-09.

[110] Benedikt Winter, Clemens Winter, Johannes Schilling, and André Bardow. A smile is all you need: predicting limiting activity coefficients from smiles with natural language processing. *Digital Discovery*, 1(6):859–869, 2022.

[111] Kevin Maik Jablonka, Philippe Schwaller, Andres Ortega-Guerrero, and Berend Smit. Leveraging large language models for predictive chemistry. *Nature Machine Intelligence*, 6(2):161–169, 2024.

[112] Niklas Gebauer, Michael Gastegger, and Kristof Schütt. Symmetry-adapted generation of 3d point sets for the targeted discovery of molecules. *Advances in neural information processing systems*, 32, 2019.

[113] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. *Advances in neural information processing systems*, 33:6840–6851, 2020.

[114] Diederik P Kingma, Max Welling, et al. Auto-encoding variational bayes, 2013.

[115] Diederik P Kingma, Max Welling, et al. An introduction to variational autoencoders. *Foundations and Trends® in Machine Learning*, 12(4):307–392, 2019.

[116] Cai-Yuan Ye, Hong-Ming Weng, and Quan-Sheng Wu. Con-cdvae: A method for the conditional generation of crystal structures. *Computational Materials Today*, 1:100003, 2024.

[117] Hang Xiao, Rong Li, Xiaoyang Shi, Yan Chen, Liangliang Zhu, Xi Chen, and Lei Wang. An invertible, invariant crystal representation for inverse design of solid-state materials using generative deep learning. *Nature Communications*, 14(1):7027, 2023.

[118] Nate Gruver, Anuroop Sriram, Andrea Madotto, Andrew Gordon Wilson, C Lawrence Zitnick, and Zachary Ulissi. Fine-tuned language models generate stable inorganic materials as text. *arXiv preprint arXiv:2402.04379*, 2024.

[119] Daniel Flam-Shepherd and Alán Aspuru-Guzik. Language models can generate molecules, materials, and protein binding sites directly in three dimensions as xyz, cif, and pdb files. *arXiv preprint arXiv:2305.05708*, 2023.

[120] Kamal Choudhary, Kevin F Garrity, Andrew CE Reid, Brian DeCost, Adam J Biacchi, Angela R Hight Walker, Zachary Trautt, Jason Hattrick-Simpers, A Gilad Kusne, Andrea Centrone, et al. Jarvis: An integrated infrastructure for data-driven materials design. *Preprint at, https://arxiv. org/abs/2007.01831*, 2020.

[121] Jingru Gan, Peichen Zhong, Yuanqi Du, Yanqiao Zhu, Chenru Duan, Haorui Wang, Carla P Gomes, Kristin A Persson, Daniel Schwalbe-Koda, and Wei Wang. Large language models are innate crystal structure generators. *arXiv preprint arXiv:2502.20933*, 2025.

[122] Philippe Schwaller, Teodoro Laino, Théophile Gaudin, Peter Bolgar, Christopher A Hunter, Costas Bekas, and Alpha A Lee. Molecular transformer: a model for uncertainty-calibrated chemical reaction prediction. *ACS central science*, 5(9):1572–1583, 2019.

[123] USPTO. Uspto – open data portal. `https://data.uspto.gov/home`, 2025.

[124] John S Delaney. Esol: estimating aqueous solubility directly from molecular structure. *Journal of chemical information and computer sciences*, 44(3):1000–1005, 2004.

[125] Zi-Yi Chen, Fan-Kai Xie, Meng Wan, Yang Yuan, Miao Liu, Zong-Guo Wang, Sheng Meng, and Yan-Gang Wang. Matchat: A large language model and application service platform for materials science. *Chinese Physics B*, 32(11):118104, 2023.

[126] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022.

[127] Chonghuan Zhang, Qianghua Lin, Biwei Zhu, Haopeng Yang, Xiao Lian, Hao Deng, Jiajun Zheng, and Kuangbiao Liao. Synask: unleashing the power of large language models in organic synthesis. *Chemical Science*, 16(1):43–56, 2025.

[128] Santiago Miret and Nandan M Krishnan. Are llms ready for real-world materials discovery? *arXiv preprint arXiv:2402.05200*, 2024.

[129] Jacob Fish, Gregory J Wagner, and Sinan Keten. Mesoscopic and multiscale modelling in materials. *Nature materials*, 20(6):774–786, 2021.

[130] Helgi I Ingólfsson, Harsh Bhatia, Fikret Aydin, Tomas Oppelstrup, Cesar A López, Liam G Stanton, Timothy S Carpenter, Sergio Wong, Francesco Di Natale, Xiaohua Zhang, et al. Machine learning-driven multiscale modeling: bridging the scales with a next-generation simulation infrastructure. *Journal of Chemical Theory and Computation*, 19(9):2658–2675, 2023.

[131] Fadi Aldakheel, Elsayed S Elsayed, Tarek I Zohdi, and Peter Wriggers. Efficient multiscale modeling of heterogeneous materials using deep neural networks. *Computational Mechanics*, 72(1):155–171, 2023.

[132] John Mayfield, Daniel Lowe, and Roger Sayle. Pistachio: Search and faceting of large reaction databases. In *ABSTRACTS OF PAPERS OF THE AMERICAN CHEMICAL SOCIETY*, volume 254. AMER CHEMICAL SOC 1155 16TH ST, NW, WASHINGTON, DC 20036 USA, 2017.

[133] Carl Edwards, ChengXiang Zhai, and Heng Ji. Text2mol: Cross-modal molecule retrieval with natural language queries. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 595–607, 2021.

[134] Jannis Born and Matteo Manica. Regression transformer enables concurrent sequence regression and generation for molecular language modelling. *Nature Machine Intelligence*, 5(4):432–444, 2023.

[135] Zhenqin Wu, Bharath Ramsundar, Evan N Feinberg, Joseph Gomes, Caleb Geniesse, Aneesh S Pappu, Karl Leswing, and Vijay Pande. Moleculenet: a benchmark for molecular machine learning. *Chemical science*, 9(2):513–530, 2018.

[136] Roshan Rao, Nicholas Bhattacharya, Neil Thomas, Yan Duan, Peter Chen, John Canny, Pieter Abbeel, and Yun Song. Evaluating protein transfer learning with tape. *Advances in neural information processing systems*, 32, 2019.

[137] Eduardo Soares, Victor Shirasuna, Emilio Vital Brazil, Renato Cerqueira, Dmitry Zubarev, and Kristin Schmidt. A large encoder-decoder family of foundation models for chemical language. *arXiv preprint arXiv:2407.20267*, 2024.

[138] Indra Priyadarsini, Seiji Takeda, Lisa Hamada, Emilio Vital Brazil, Eduardo Soares, and Hajime Shinohara. Self-bart: A transformer-based molecular representation model using selfies. *arXiv preprint arXiv:2410.12348*, 2024.

[139] Akihiro Kishimoto, Hiroshi Kajino, Masataka Hirose, Junta Fuchiwaki, Indra Priyadarsini, Lisa Hamada, Hajime Shinohara, Daiju Nakano, and Seiji Takeda. Mhg-gnn: Combination of molecular hypergraph grammar with graph neural network. *arXiv preprint arXiv:2309.16374*, 2023.

[140] Andrew Trotman, Antti Puurula, and Blake Burgess. Improvements to bm25 and language models examined. In *Proceedings of the 19th Australasian Document Computing Symposium*, pages 58–65, 2014.

[141] Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. Unsupervised dense information retrieval with contrastive learning. *arXiv preprint arXiv:2112.09118*, 2021.

[142] Victor Fung, Jiaxin Zhang, Eric Juarez, and Bobby G Sumpter. Benchmarking graph neural networks for materials chemistry. *npj Computational Materials*, 7(1):84, 2021.

[143] Philipp Thölke and Gianni De Fabritiis. Torchmd-net: equivariant transformers for neural network based molecular potentials. *arXiv preprint arXiv:2202.02541*, 2022.

[144] Dejan Zagorac, H Müller, S Ruehl, J Zagorac, and Silke Rehme. Recent developments in the inorganic crystal structure database: theoretical crystal structure data and related features. *Applied Crystallography*, 52(5): 918–925, 2019.

[145] Scott Kirklin, James E Saal, Bryce Meredig, Alex Thompson, Jeff W Doak, Muratahan Aykol, Stephan Rühl, and Chris Wolverton. The open quantum materials database (oqmd): assessing the accuracy of dft formation energies. *npj Computational Materials*, 1(1):1–15, 2015.

[146] James E Saal, Scott Kirklin, Muratahan Aykol, Bryce Meredig, and Christopher Wolverton. Materials design and discovery with high-throughput density functional theory: the open quantum materials database (oqmd). *Jom*, 65:1501–1509, 2013.

[147] Claudia Draxl and Matthias Scheffler. The nomad laboratory: from data sharing to artificial intelligence. *Journal of Physics: Materials*, 2(3):036001, 2019.

[148] Luis Barroso-Luque, Muhammed Shuaibi, Xiang Fu, Brandon M Wood, Misko Dzamba, Meng Gao, Ammar Rizvi, C Lawrence Zitnick, and Zachary W Ulissi. Open materials 2024 (omat24) inorganic materials dataset and models. *arXiv preprint arXiv:2410.12771*, 2024.

[149] SNU MDIL. Snumat. `https://www.snumat.com/apis`, 2025.

[150] Bowen Deng, Peichen Zhong, KyuJung Jun, Janosh Riebesell, Kevin Han, Christopher J Bartel, and Gerbrand Ceder. Chgnet as a pretrained universal neural network potential for charge-informed atomistic modelling. *Nature Machine Intelligence*, 5(9):1031–1041, 2023.

[151] Jonathan Schmidt, Tiago FT Cerqueira, Aldo H Romero, Antoine Loew, Fabian Jäger, Hai-Chen Wang, Silvana Botti, and Miguel AL Marques. Improving machine-learning models in materials science through large datasets. *Materials Today Physics*, 48:101560, 2024.

[152] Nathan Brown, Marco Fiscato, Marwin HS Segler, and Alain C Vaucher. Guacamol: benchmarking models for de novo molecular design. *Journal of chemical information and modeling*, 59(3):1096–1108, 2019.

[153] Daniil Polykovskiy, Alexander Zhebrak, Benjamin Sanchez-Lengeling, Sergey Golovanov, Oktai Tatanov, Stanislav Belyaev, Rauf Kurbanov, Aleksey Artamonov, Vladimir Aladinskiy, Mark Veselov, et al. Molecular sets (moses): a benchmarking platform for molecular generation models. *Frontiers in pharmacology*, 11:565644, 2020.

[154] Vahe Tshitoyan, John Dagdelen, Leigh Weston, Alexander Dunn, Ziqin Rong, Olga Kononova, Kristin A Persson, Gerbrand Ceder, and Anubhav Jain. Unsupervised word embeddings capture latent knowledge from materials science literature. *Nature*, 571(7763):95–98, 2019.

[155] Leigh Weston, Vahe Tshitoyan, John Dagdelen, Olga Kononova, Amalie Trewartha, Kristin A Persson, Gerbrand Ceder, and Anubhav Jain. Named entity recognition and normalization applied to large-scale information extraction from the materials science literature. *Journal of chemical information and modeling*, 59(9):3692–3702, 2019.

[156] Yu Song, Santiago Miret, and Bang Liu. Matsci-nlp: Evaluating scientific language models on materials science language tasks using text-to-schema modeling. *arXiv preprint arXiv:2305.08264*, 2023.

[157] Andre Niyongabo Rubungo, Kangming Li, Jason Hattrick-Simpers, and Adji Bousso Dieng. Llm4mat-bench: Benchmarking large language models for materials property prediction. *Machine Learning: Science and Technology*, 2024.

[158] Zhiqiang Zhong, Kuangyu Zhou, and Davide Mottin. Benchmarking large language models for molecule prediction tasks. *arXiv preprint arXiv:2403.05075*, 2024.

[159] Nawaf Alampara, Indrajeet Mandal, Pranav Khetarpal, Hargun Singh Grover, Mara Schilling-Wilhelmi, N M Anoop Krishnan, and Kevin Maik Jablonka. MaCBench: A multimodal chemistry and materials science benchmark. In *AI for Accelerated Materials Design - NeurIPS 2024*, 2024. URL `https://openreview.net/forum?id=Q2PNocDcp6`.

[160] Maurice Weber, Dan Fu, Quentin Anthony, Yonatan Oren, Shane Adams, Anton Alexandrov, Xiaozhong Lyu, Huu Nguyen, Xiaozhe Yao, Virginia Adams, et al. Redpajama: an open dataset for training large language models. *Advances in neural information processing systems*, 37:116462–116492, 2024.

[161] Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *Proceedings of the IEEE international conference on computer vision*, pages 19–27, 2015.

[162] Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. *arXiv preprint arXiv:1609.07843*, 2016.

[163] Common Crawl Team. Common crawl dataset, 2025. URL `https://commoncrawl.org/`.

[164] Arman Cohan, Sergey Feldman, Iz Beltagy, Doug Downey, and Daniel S Weld. Specter: Document-level representation learning using citation-informed transformers. *arXiv preprint arXiv:2004.07180*, 2020.

[165] Ting-Yao Hsu, C Lee Giles, and Ting-Hao'Kenneth' Huang. Scicap: Generating captions for scientific figures. *arXiv preprint arXiv:2110.11624*, 2021.

[166] Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. Orca: Progressive learning from complex explanation traces of gpt-4. *arXiv preprint arXiv:2306.02707*, 2023.

[167] Xiang Yue, Tianyu Zheng, Ge Zhang, and Wenhu Chen. Mammoth2: Scaling instructions from the web. *Advances in Neural Information Processing Systems*, 37:90629–90660, 2024.

[168] Aida Amini, Saadia Gabriel, Peter Lin, Rik Koncel-Kedziorski, Yejin Choi, and Hannaneh Hajishirzi. Mathqa: Towards interpretable math word problem solving with operation-based formalisms. *arXiv preprint arXiv:1905.13319*, 2019.

[169] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint arXiv:1803.05457*, 2018.

[170] Yonatan Bisk, Rowan Zellers, Jianfeng Gao, Yejin Choi, et al. Piqa: Reasoning about physical commonsense in natural language. In *Proceedings of the AAAI conference on artificial intelligence*, volume 34, pages 7432–7439, 2020.

[171] Johannes Welbl, Nelson F Liu, and Matt Gardner. Crowdsourcing multiple choice science questions. *arXiv preprint arXiv:1707.06209*, 2017.

[172] Pan Lu, Swaroop Mishra, Tanglin Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. Learn to explain: Multimodal reasoning via thought chains for science question answering. *Advances in Neural Information Processing Systems*, 35:2507–2521, 2022.

[173] Casper W Andersen, Rickard Armiento, Evgeny Blokhin, Gareth J Conduit, Shyam Dwaraknath, Matthew L Evans, Ádám Fekete, Abhijith Gopakumar, Saulius Gražulis, Andrius Merkys, et al. Optimade, an api for exchanging materials data. *Scientific data*, 8(1):217, 2021.

[174] Shyue Ping Ong, William Davidson Richards, Anubhav Jain, Geoffroy Hautier, Michael Kocher, Shreyas Cholia, Dan Gunter, Vincent L Chevrier, Kristin A Persson, and Gerbrand Ceder. Python materials genomics (pymatgen): A robust, open-source python library for materials analysis. *Computational Materials Science*, 68:314–319, 2013.

[175] Yuanqi Du, Yingheng Wang, Yining Huang, Jianan Canal Li, Yanqiao Zhu, Tian Xie, Chenru Duan, John Gregoire, and Carla P Gomes. M$^2$hub: Unlocking the potential of machine learning for materials discovery. *Advances in Neural Information Processing Systems*, 36:77359–77378, 2023.

[176] Anubhav Jain, Shyue Ping Ong, Wei Chen, Bharat Medasani, Xiaohui Qu, Michael Kocher, Miriam Brafman, Guido Petretto, Gian-Marco Rignanese, Geoffroy Hautier, et al. Fireworks: a dynamic workflow system designed for high-throughput applications. *Concurrency and Computation: Practice and Experience*, 27(17):5037–5059, 2015.

[177] The Materials Project. Maggma Toolkit. `https://materialsproject.github.io/maggma/`, 2025.

[178] Alexander Dunn, Qi Wang, Alex Ganose, Daniel Dopp, and Anubhav Jain. Benchmarking materials property prediction methods: the matbench test set and automatminer reference algorithm. *npj Computational Materials*, 6(1):138, 2020.

[179] Kamal Choudhary, Brian DeCost, Lily Major, Keith Butler, Jeyan Thiyagalingam, and Francesca Tavazza. Unified graph neural network force-field for the periodic table: solid state applications. *Digital Discovery*, 2(2):346–355, 2023.

[180] Ollama. Ollama tool. `https://ollama.com/`, 2025.

[181] LangChain. Langchain tool. `https://github.com/langchain-ai/langchain`, 2025.

[182] Matthew Horton, Jimmy-Xuan Shen, Jordan Burns, Orion Cohen, François Chabbey, Alex M Ganose, Rishabh Guha, Patrick Huck, Hamming Howard Li, Matthew McDermott, et al. Crystal toolkit: A web app framework to improve usability and accessibility of materials science research algorithms. *arXiv preprint arXiv:2302.06147*, 2023.

[183] Kiran Mathew, Joseph H Montoya, Alireza Faghaninia, Shyam Dwarakanath, Muratahan Aykol, Hanmei Tang, Iek-heng Chu, Tess Smidt, Brandon Bocklund, Matthew Horton, et al. Atomate: A high-level interface to generate, execute, and analyze computational materials science workflows. *Computational Materials Science*, 139:140–152, 2017.

[184] Andrew S Rosen, Max Gallant, Janine George, Janosh Riebesell, Hrushikesh Sahasrabuddhe, Jimmy-Xuan Shen, Mingjian Wen, Matthew L Evans, Guido Petretto, David Waroquiers, et al. Jobflow: Computational workflows made simple. *Journal of Open Source Software*, 9(93):5995, 2024.

[185] The Materials Project. Emmet Toolkit. `https://materialsproject.github.io/emmet/`, 2025.

[186] Microsoft. Autogen tool. `https://github.com/microsoft/autogen`, 2025.

[187] CrewAI. Crewai tool. `https://github.com/crewAIInc/crewAI`, 2025.

[188] LlamaIndex. Llamaindex tool. `https://github.com/run-llama/llama_index`, 2025.