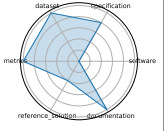
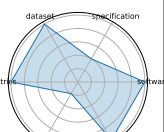

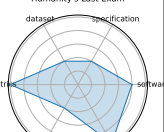
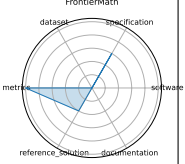
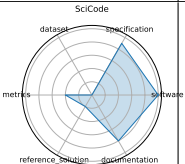
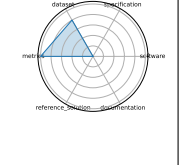
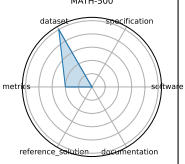


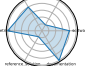
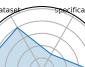
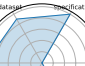
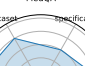
1 Benchmark Overview Table

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	MMLU (Massive Multitask Language Understanding)	Multidomain	Academic knowledge and reasoning across 57 subjects	multitask, multiple-choice, zero-shot, few-shot, knowledge probing	Multiple choice	General reasoning, subject-matter understanding	Accuracy	GPT-4o, Gemini 1.5 Pro, o1, DeepSeek-R1	[1]⇒
	GPQA Diamond	Science	Graduate-level scientific reasoning	Google-proof, graduate-level, science QA, chemistry, physics	Multiple choice, Multi-step QA	Scientific reasoning, deep knowledge	Accuracy	o1, DeepSeek-R1	[2]⇒
	ARC-Challenge (Advanced Reasoning Challenge)	Science	Grade-school science with reasoning emphasis	grade-school, science QA, challenge set, reasoning	Multiple choice	Commonsense and scientific reasoning	Accuracy	GPT-4, Claude	[3]⇒
	Humanity's Last Exam	Multidomain	Broad cross-domain academic reasoning	cross-domain, academic exam, multiple-choice, multi-disciplinary	Multiple choice	Cross-domain academic reasoning	Accuracy	unknown	[4]⇒

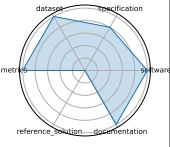
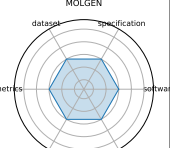
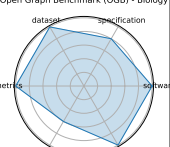
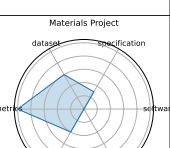
Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	FrontierMath	Mathematics	Challenging advanced mathematical reasoning	symbolic reasoning, number theory, algebraic geometry, category theory	Problem solving	Symbolic and abstract mathematical reasoning	Accuracy	unkown	[5]⇒
	SciCode	Scientific Programming	Scientific code generation and problem solving	code synthesis, scientific computing, programming benchmark	Coding	Program synthesis, scientific computing	Solve rate (%)	Claude3.5-Sonnet	[6]⇒
	AIME (American Invitational Mathematics Examination)	Mathematics	Pre-college advanced problem solving	algebra, combinatorics, number theory, geometry	Problem solving	Mathematical problem-solving and reasoning	Accuracy	unkown	[7]⇒
	MATH-500	Mathematics	Math reasoning generalization	calculus, algebra, number theory, geometry	Problem solving	Math reasoning and generalization	Accuracy	unkown	[8]⇒

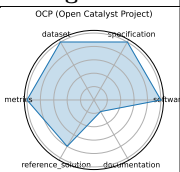
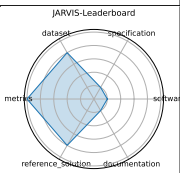
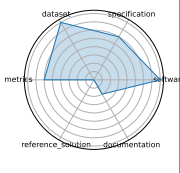
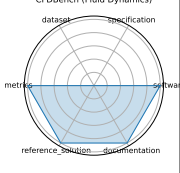
Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	CURIE (Scientific Long-Context Understanding, Reasoning and Information Extraction)	Multidomain Science	Long-context scientific reasoning	long-context, information extraction, multimodal	Information extraction, Reasoning, Concept tracking, Aggregation, Algebraic manipulation, Multimodal comprehension	Long-context understanding and scientific reasoning	Accuracy	unknown	[9]⇒
	FEABench (Finite Element Analysis Benchmark)	Computational Engineering	FEA simulation accuracy and performance	finite element, simulation, PDE	Simulation, Performance evaluation	Numerical simulation accuracy and efficiency	Solve time, Error norm	FEniCS, deal.II	[10]⇒
	SPIQA (Scientific Paper Image Question Answering)	Computer Science	Multimodal QA on scientific figures	multimodal QA, figure understanding, table comprehension, chain-of-thought	Question answering, Multimodal QA, Chain-of-Thought evaluation	Visual-textual reasoning in scientific contexts	Accuracy, F1 score	Chain-of-Thought models, Multimodal QA systems	[11]⇒
	MedQA	Medical Question Answering	Medical board exam QA	USMLE, diagnostic QA, medical knowledge, multilingual	Multiple choice	Medical diagnosis and knowledge retrieval	Accuracy	Neural reader, Retrieval-based QA systems	[12]⇒

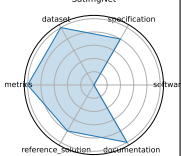
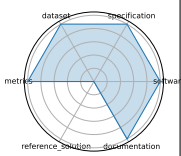
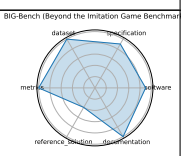
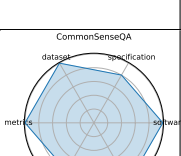
Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	BaisBench (Biological AI Scientist Benchmark)	Computational Biology	Omics- driven AI research tasks	single-cell annotation, biological QA, au- tonomous discovery	Cell type anno- tation, Multiple choice	Autonomous bi- ological research capabilities	Annotation accuracy, QA accu- racy	LLM-based AI scientist agents	[13]⇒
	MOLGEN	Computational Chemistry	Molecular generation and opti- mization	SELFIES, GAN, prop- erty opti- mization	Distribution learning, Goal- oriented genera- tion	Generation of valid and opti- mized molecular structures	Validity%, Novelty%, QED, Docking score	MolGen	[14]⇒
	Open Graph Benchmark (OGB) - Biology	Graph ML	Biological graph property prediction	node predic- tion, link pre- diction, graph classification	Node prop- erty prediction, Link property prediction, Graph property prediction	Scalability and generalization in graph ML for bi- ology	Accuracy, ROC-AUC	GCN, Graph- SAGE, GAT	[15]⇒
	Materials Project	Materials Science	DFT-based property prediction	DFT, ma- terials genome, high- throughput	Property predic- tion	Prediction of in- organic material properties	MAE, R ²	Automatminer Crystal Graph Neural Networks	[16]⇒

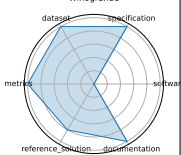
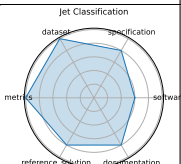
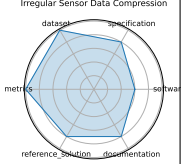
Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	OCP (Open Catalyst Project)	Chemistry; Materials Science	Catalyst adsorption energy prediction	DFT relaxations, adsorption energy, graph neural networks	Energy prediction, Force prediction	Prediction of adsorption energies and forces	MAE (energy), MAE (force)	CGCNN, SchNet, DimeNet++, GemNet-OC	[17]–[20]⇒
	JARVIS-Leaderboard	Materials Science; Benchmarking	Comparative evaluation of materials design methods	leaderboards, materials methods, simulation	Method benchmarking, Leaderboard ranking	Performance comparison across diverse materials design methods	MAE, RMSE, Accuracy	unkown	[21]⇒
	Quantum Computing Benchmarks (QML)	Quantum Computing	Quantum algorithm performance evaluation	quantum circuits, state preparation, error correction	Circuit benchmarking, State classification	Quantum algorithm performance and fidelity	Fidelity, Success probability	IBM Q, IonQ, AQT@LBNL	[22]⇒
	CFDBench (Fluid Dynamics)	Fluid Dynamics; Scientific ML	Neural operator surrogate modeling	neural operators, CFD, FNO, DeepONet	Surrogate modeling	Generalization of neural operators for PDEs	L2 error, MAE	FNO, DeepONet, U-Net	[23]⇒

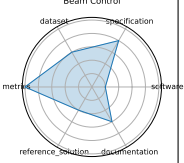
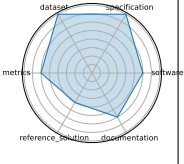
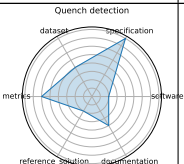
Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	SatImgNet	Remote Sensing	Satellite imagery classification	land-use, zero-shot, multi-task	Image classification	Zero-shot land-use classification	Accuracy	CLIP, BLIP, ALBEF	[24]⇒
	ClimateLearn	Climate Science; Forecasting	ML for weather and climate modeling	medium-range forecasting, ERA5, data-driven	Forecasting	Global weather prediction (3-5 days)	RMSE, Anomaly correlation	CNN baselines, ResNet variants	[25]⇒
	BIG-Bench (Beyond the Imitation Game Benchmark)	NLP; AI Evaluation	Diverse reasoning and generalization tasks	few-shot, multi-task, bias analysis	Few-shot evaluation, Multi-task evaluation	Reasoning and generalization across diverse tasks	Accuracy, Task-specific metrics	GPT-3, Dense Transformers, Sparse Transformers	[26]⇒
	CommonSenseQA	NLP; Commonsense	Commonsense question answering	Commonsense ConceptNet, multiple-choice, adversarial	Multiple choice	Commonsense reasoning and knowledge integration	Accuracy	BERT-large, RoBERTa, GPT-3	[27]⇒


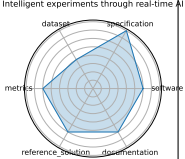
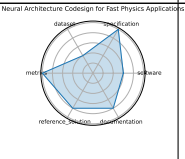
Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	Winogrande	NLP; Commonsense	Winograd Schema-style pronoun resolution	adversarial, pronoun resolution	Pronoun resolution	Robust commonsense reasoning	Accuracy, AUC	RoBERTa, BERT, GPT-2	[28]⇒
	Jet Classification	Particle Physics	Real-time classification of particle jets using HL-LHC simulation features	classification, real-time ML, jet tagging, QKeras	Classification	Real-time inference, model compression performance	Accuracy, AUC	Keras DNN, QKeras quantized DNN	[29]⇒
	Irregular Sensor Data Compression	Particle Physics	Real-time compression of sparse sensor data with autoencoders	compression, autoencoder, sparse data, irregular sampling	Compression	Reconstruction quality, compression efficiency	MSE, Compression ratio	Autoencoder, Quantized autoencoder	[30]⇒

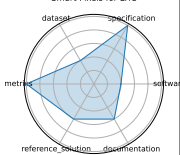
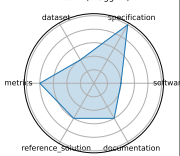
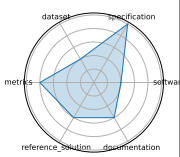
Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	Beam Control	Accelerators and Magnets	Reinforcement learning control of accelerator beam position	RL, beam stabilization, control systems, simulation	Control	Policy performance in simulated accelerator control	Stability, Control loss	DDPG, PPO (planned)	[31], [32]⇒
	Ultrafast jet classification at the HL-LHC	Particle Physics	FPGA-optimized real-time jet origin classification at the HL-LHC	jet classification, FPGA, quantization-aware training, Deep Sets, Interaction Networks	Classification	Real-time inference under FPGA constraints	Accuracy, Latency, Resource utilization	MLP, Deep Sets, Interaction Network	[33]⇒
	Quench detection	Accelerators and Magnets	Real-time detection of superconducting magnet quenches using ML	quench detection, autoencoder, anomaly detection, real-time	Anomaly detection, Quench localization	Real-time anomaly detection with multi-modal sensors	ROC-AUC, Detection latency	Autoencoder, RL agents (in development)	[34]⇒

Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	DUNE	Particle Physics	Real-time ML for DUNE DAQ time-series data	DUNE, time-series, real-time, trigger	Trigger selection, Time-series anomaly detection	Low-latency event detection	Detection efficiency, Latency	CNN, LSTM (planned)	[35]⇒
	Intelligent experiments through real-time AI	Instrumentation and Detectors; Nuclear Physics; Particle Physics	Real-time FPGA-based triggering and detector control for sPHENIX and future EIC	FPGA, Graph Neural Network, hls4ml, real-time inference, detector control	Trigger classification, Detector control, Real-time inference	Low-latency GNN inference on FPGA	Accuracy (charm and beauty detection), Latency (micros), Resource utilization (LUT/FF/BRAM/DSP)	Bipartite Graph Network with Set Transformers (BGN-ST), GarNet (edge-AE/ISP)	[36]⇒
	Neural Architecture Codesign for Fast Physics Applications	Physics; Materials Science; Particle Physics	Automated neural architecture search and hardware-efficient model codesign for fast physics applications	neural architecture search, FPGA deployment, quantization, pruning, hls4ml	Classification, Peak finding	Hardware-aware model optimization; low-latency inference	Accuracy, Latency, Resource utilization	NAC-based BraggNN, NAC-optimized Deep Sets (jet)	[37]⇒

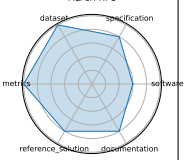
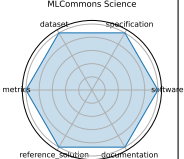
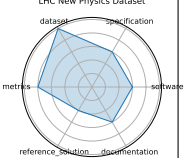
Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	Smart Pixels for LHC	Particle Physics; Instrumentation and Detectors	On-sensor, in-pixel ML filtering for high-rate LHC pixel detectors	smart pixel, on-sensor inference, data reduction, trigger	Image Classification, Data filtering	On-chip, low-power inference; data reduction	Data rejection rate, Power per pixel	2-layer pixel NN	[38]⇒
	HEDM (BraggNN)	Material Science	Fast Bragg peak analysis using deep learning in diffraction microscopy	BraggNN, diffraction, peak finding, HEDM	Peak detection	High-throughput peak localization	Localization accuracy, Inference time	BraggNN	[39]⇒
	4D-STEM	Material Science	Real-time ML for scanning transmission electron microscopy	4D-STEM, electron microscopy, real-time, image processing	Image Classification, Streamed data inference	Real-time large-scale microscopy inference	Classification accuracy, Throughput	CNN models (prototype)	[40]⇒

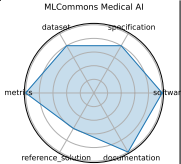
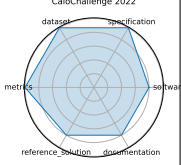
Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	In-Situ High-Speed Computer Vision	Fusion/Plasma	Real-time image classification for in-situ plasma diagnostics	plasma, in-situ vision, real-time ML	Image Classification	Real-time diagnostic inference	Accuracy, FPS	CNN	[41]⇒
	BenchCouncil AIBench	General	End-to-end AI benchmarking across micro, component, and application levels	benchmarking, AI systems, application-level evaluation	Training, Inference, End-to-end AI workloads	System-level AI workload performance	Throughput, Latency, Accuracy	ResNet, BERT, GANs, Recommendation systems	[42]⇒
	BenchCouncil Big-DataBench	General	Big data and AI benchmarking across structured, semi-structured, and unstructured data workloads	big data, AI benchmarking, data analytics	Data pre-processing, Inference, End-to-end data pipelines	Data processing and AI model inference performance at scale	Data throughput, Latency, Accuracy	CNN, LSTM, SVM, XGBoost	[43]⇒

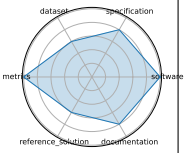
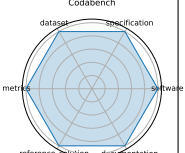
Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	MLPerf HPC	Cosmology, Climate, Protein Structure, Catalysis	Scientific ML training and inference on HPC systems	HPC, training, inference, scientific ML	Training, Inference	Scaling efficiency, training time, model accuracy on HPC	Training time, Accuracy, GPU utilization	CosmoFlow, DeepCAM, OpenCatalyst	[44]⇒
	MLCommons Science	Earthquake, Satellite Image, Drug Discovery, Electron Microscope, CFD	AI benchmarks for scientific applications including time-series, imaging, and simulation	science AI, benchmark, MLCommons, HPC	Time-series analysis, Image classification, Simulation surrogate modeling	Inference accuracy, simulation speed-up, generalization	MAE, Accuracy, Speedup vs simulation	CNN, GNN, Transformer	[45]⇒
	LHC New Physics Dataset	Particle Physics; Real-time Triggering	Real-time LHC event filtering for anomaly detection using proton collision data	anomaly detection, proton collision, real-time inference, event filtering, unsupervised ML	Anomaly detection, Event classification	Unsupervised signal detection under latency and bandwidth constraints	ROC-AUC, Detection efficiency	Autoencoder, Variational autoencoder, Isolation forest	[46]⇒

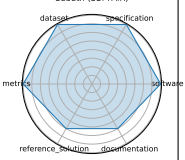
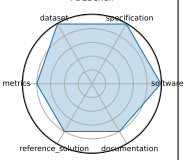
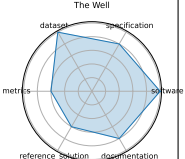
Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	MLCommons Medical AI	Healthcare; Medical AI	Federated benchmarking and evaluation of medical AI models across diverse real-world clinical data	medical AI, federated evaluation, privacy-preserving, fairness, healthcare benchmarks	Federated evaluation, Model validation	Clinical accuracy, fairness, generalizability, privacy compliance	ROC AUC, Accuracy, Fairness metrics	MedPerf-validated CNNs, GaNDLF workflows	[47]⇒
	CaloChallenge 2022	LHC Calorimeter; Particle Physics	Fast generative-model-based calorimeter shower simulation evaluation	calorimeter simulation, generative models, surrogate modeling, LHC, fast simulation	Surrogate modeling	Simulation fidelity, speed, efficiency	Histogram similarity, Classifier AUC, Generation latency	VAE variants, GAN variants, Normalizing flows, Diffusion models	[48]⇒

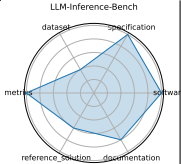
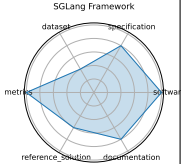
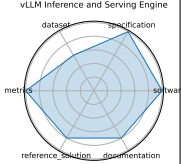
Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
 <p>Papers With Code (SOTA Platform)</p>	Papers With Code (SOTA Platform)	General ML; All domains	Open platform tracking state-of-the-art results, benchmarks, and implementations across ML tasks and papers	leaderboard, benchmarking, reproducibility, open-source	Multiple (Classification, Detection, NLP, etc.)	Model performance across tasks (accuracy, F1, BLEU, etc.)	Task-specific (Accuracy, F1, BLEU, etc.)	All published models with code	[49]⇒
 <p>Codabench</p>	Codabench	General ML; Multiple	Open-source platform for organizing reproducible AI benchmarks and competitions	benchmark platform, code submission, competitions, meta-benchmark	Multiple	Model reproducibility, performance across datasets	Submission count, Leaderboard ranking, Task-specific metrics	Arbitrary code submissions	[50]⇒

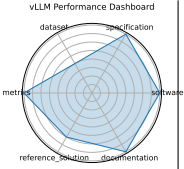
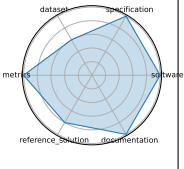
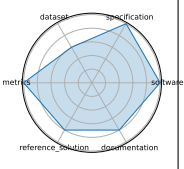
Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	Sabath (SBI-FAIR)	Systems; Metadata	FAIR metadata framework for ML-driven surrogate workflows in HPC systems	meta-benchmark, metadata, HPC, surrogate modeling	Systems benchmarking	Metadata tracking, reproducible HPC workflows	Metadata completeness, FAIR compliance	NA	[51]⇒
	PDEBench	CFD; Weather Modeling	Benchmark suite for ML-based surrogates solving time-dependent PDEs	PDEs, CFD, scientific ML, surrogate modeling, NeurIPS	Supervised Learning	Time-dependent PDE modeling; physical accuracy	RMSE, boundary RMSE, Fourier RMSE	FNO, U-Net, PINN, Gradient-Based inverse methods	[52]⇒
	The Well	biological systems, fluid dynamics, acoustic scattering, astrophysical MHD	Foundation model + surrogate dataset spanning 16 physical simulation domains	surrogate modeling, foundation model, physics simulations, spatiotemporal dynamics	Supervised Learning	Surrogate modeling, physics-based prediction	Dataset size, Domain breadth	FNO baselines, U-Net baselines	[53]⇒

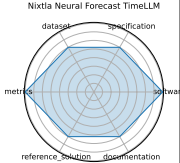
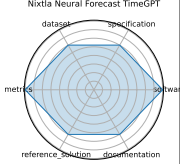
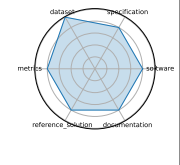
Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	LLM-Inference-Bench	LLM; HPC/inference	Hardware performance benchmarking of LLMs on AI accelerators	LLM, inference benchmarking, GPU, accelerator, throughput	Inference Benchmarking	Inference throughput, latency, hardware utilization	Token throughput (tok/s), Latency, Framework-hardware mix performance	LLaMA-2-7B, LLaMA-2-70B, Mistral-7B, Qwen-7B	[54]⇒
	SGLang Framework	LLM Vision	Fast serving framework for LLMs and vision-language models	LLM serving, vision-language, RadixAttention, performance, JSON decoding	Model serving framework	Serving throughput, JSON/task-specific latency	Tokens/sec, Time-to-first-token, Throughput gain vs baseline	LLaVA, DeepSeek, Llama	[55]⇒
	vLLM Inference and Serving Engine	LLM; HPC/inference	High-throughput, memory-efficient inference and serving engine for LLMs	LLM inference, PagedAttention, CUDA graph, streaming API, quantization	Inference Benchmarking	Throughput, latency, memory efficiency	Tokens/sec, Time to First Token (TTFT), Memory footprint	LLaMA, Mixtral, FlashAttention-based models	[56]⇒

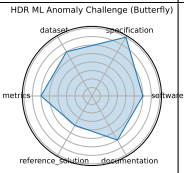
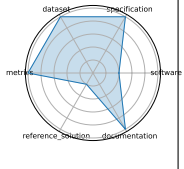
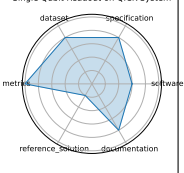
Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	vLLM Performance Dashboard	LLM; HPC/inference	Interactive dashboard showing inference performance of vLLM	Dashboard, Throughput visualization, Latency analysis, Metric tracking	Performance visualization	Throughput, latency, hardware utilization	Tokens/sec, TTFT, Memory usage	LLaMA-2, Mistral, Qwen	[57]⇒
	Nixtla NeuralForecast	Time-series forecasting; General ML	High-performance neural forecasting library with >30 models	time-series, neural forecasting, NBEATS, NHITS, TFT, probabilistic forecasting, usability	Time-series forecasting	Forecast accuracy, interpretability, speed	RMSE, MAPE, CRPS	NBEATS, NHITS, TFT, DeepAR	[58]⇒
	Nixtla Neural Forecast NHITS	Time-series; General ML	Official NHITS implementation for long-horizon time series forecasting	NHITS, long-horizon forecasting, neural interpolation, time-series	Time-series forecasting	Accuracy, compute efficiency for long series	RMSE, MAPE	NHITS	[59]⇒


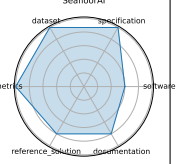
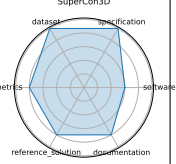
Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	Nixtla Neural Forecast TimeLLM	Time-series; General ML	Reprogramming LLMs for time series forecasting	Time-LLM, language model, time-series, reprogramming	Time-series forecasting	Model reuse via LLM, few-shot forecasting	RMSE, MAPE	Time-LLM	[60]⇒
	Nixtla Neural Forecast TimeGPT	Time-series; General ML	Time-series foundation model "TimeGPT" for forecasting and anomaly detection	TimeGPT, foundation model, time-series, generative model	Time-series forecasting, Anomaly detection	Zero-shot forecasting, anomaly detection	RMSE, Anomaly detection metrics	TimeGPT	[61]⇒
	HDR ML Anomaly Challenge (Gravitational Waves)	Astrophysics; Time-series	Detecting anomalous gravitational wave signals from LIGO/Virgo datasets	anomaly detection, gravitational waves, astrophysics, time-series	Anomaly detection	Novel event detection in physical signals	ROC-AUC, Precision/Recall	Deep latent CNNs, Autoencoders	[62]⇒

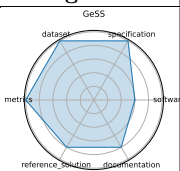
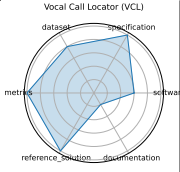
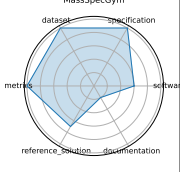
Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	HDR ML Anomaly Challenge (Butterfly)	Genomics; Image/CV	Detecting hybrid butterflies via image anomaly detection in genomic-informed dataset	anomaly detection, computer vision, genomics, butterfly hybrids	Anomaly detection	Hybrid detection in biological systems	Classification accuracy, F1 score	CNN-based detectors	[63]⇒
	HDR ML Anomaly Challenge (Sea Level Rise)	Climate Science; Time-series, Image/CV	Detecting anomalous sea-level rise and flooding events via time-series and satellite imagery	anomaly detection, climate science, sea-level rise, time-series, remote sensing	Anomaly detection	Detection of environmental anomalies	ROC-AUC, Precision/Recall	CNNs, RNNs, Transformers	[64]⇒
	Single Qubit Readout on QICK System	Quantum Computing	Real-time single-qubit state classification using FPGA firmware	qubit readout, hls4ml, FPGA, QICK	Classification	Single-shot fidelity, inference latency	Accuracy, Latency	hls4ml quantized NN	[65]⇒

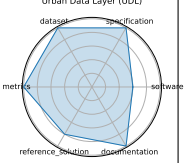
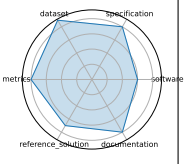

Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	GPQA: A Graduate-Level Google-Proof Question and Answer Benchmark	Science (Biology, Physics, Chemistry)	Graduate-level, expert-validated multiple-choice questions hard even with web access	Google-proof, multiple-choice, expert reasoning, science QA	Multiple choice	Scientific reasoning, knowledge probing	Accuracy	GPT-4 baseline	[66]⇒
	SeaFloorAI	Marine Science; Vision-Language	Large-scale vision-language dataset for seafloor mapping and geological classification	sonar imagery, vision-language, seafloor mapping, segmentation, QA	Image segmentation, Vision-language QA	Geospatial understanding, multimodal reasoning	Segmentation pixel accuracy, QA accuracy	SegFormer, ViLT-style multi-modal models	[67]⇒
	SuperCon3D	Materials Science; Superconductivity	Dataset and models for predicting and generating high-Tc superconductors using 3D crystal structures	superconductivity, crystal structures, equivariant GNN, generative models	Regression (Tc prediction), Generative modeling	Structure-to-property prediction, structure generation	MAE (Tc), Validity of generated structures	SODNet, DiffCSP-SC	[68]⇒

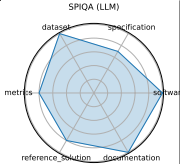
Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	GeSS	Scientific ML; Geometric Deep Learning	Benchmark suite evaluating geometric deep learning models under real-world distribution shifts	geometric deep learning, distribution shift, OOD robustness, scientific applications	Classification, Regression	OOD performance in scientific settings	Accuracy, RMSE, OOD robustness delta	GCN, EGNN, DimeNet++	[69]⇒
	Vocal Call Locator (VCL)	Neuroscience; Bioacoustics	Benchmarking sound-source localization of rodent vocalizations from multi-channel audio	source localization, bioacoustics, time-series, SSL	Sound source localization	Source localization accuracy in bioacoustic settings	Localization error (cm), Recall/Precision	CNN-based SSL models	[70]⇒
	MassSpecGym	Cheminformatics; Molecular Discovery	Benchmark suite for discovery and identification of molecules via MS/MS	mass spectrometry, molecular structure, de novo generation, retrieval, dataset	De novo generation, Retrieval, Simulation	Molecular identification and generation from spectral data	Structure accuracy, Retrieval precision, Simulation MSE	Graph-based generative models, Retrieval baselines	[71]⇒

Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	Urban Data Layer (UDL)	Urban Computing; Data Engineering	Unified data pipeline for multi-modal urban science research	data pipeline, urban science, multi-modal, benchmark	Prediction, Classification	Multi-modal urban inference, standardization	Task-specific accuracy or RMSE	Baseline regression/classification pipelines	[72]⇒
	Delta Squared-DFT	Computational Chemistry; Materials Science	Benchmarking density functional theory, Delta Squared-ML correction, reaction energetics, quantum chemistry	machine-learning corrections to DFT using Delta Squared-trained models for reaction energies	Regression	High-accuracy energy prediction, DFT correction	Mean Absolute Error (eV), Energy ranking accuracy	Delta Squared-ML correction networks, Kernel ridge regression	[73]⇒
	LLMs for Crop Science	Agricultural Science; NLP	Evaluating LLMs on crop trait QA and textual inference tasks with domain-specific prompts	crop science, prompt engineering, domain adaptation, question answering	Question Answering, Inference	Scientific knowledge, crop reasoning	Accuracy, F1 score	GPT-4, LLaMA-2-13B, T5-XXL	[74]⇒

Continued on next page

Ratings	Name	Domain	Focus	Keywords	Task Types	AI Capability	Metrics	Models	Citation
	SPIQA (LLM)	Multimodal Scientific QA; Computer Vision	Evaluating LLMs on image-based scientific paper figure QA tasks (LLM Adapter performance)	multimodal QA, scientific figures, image+text, chain-of-thought prompting	Multimodal QA	Visual reasoning, scientific figure understanding	Accuracy, F1 score	LLaVA, MiniGPT-4, Owl-LLM adapter variants	[75]⇒

2 Radar Chart Table

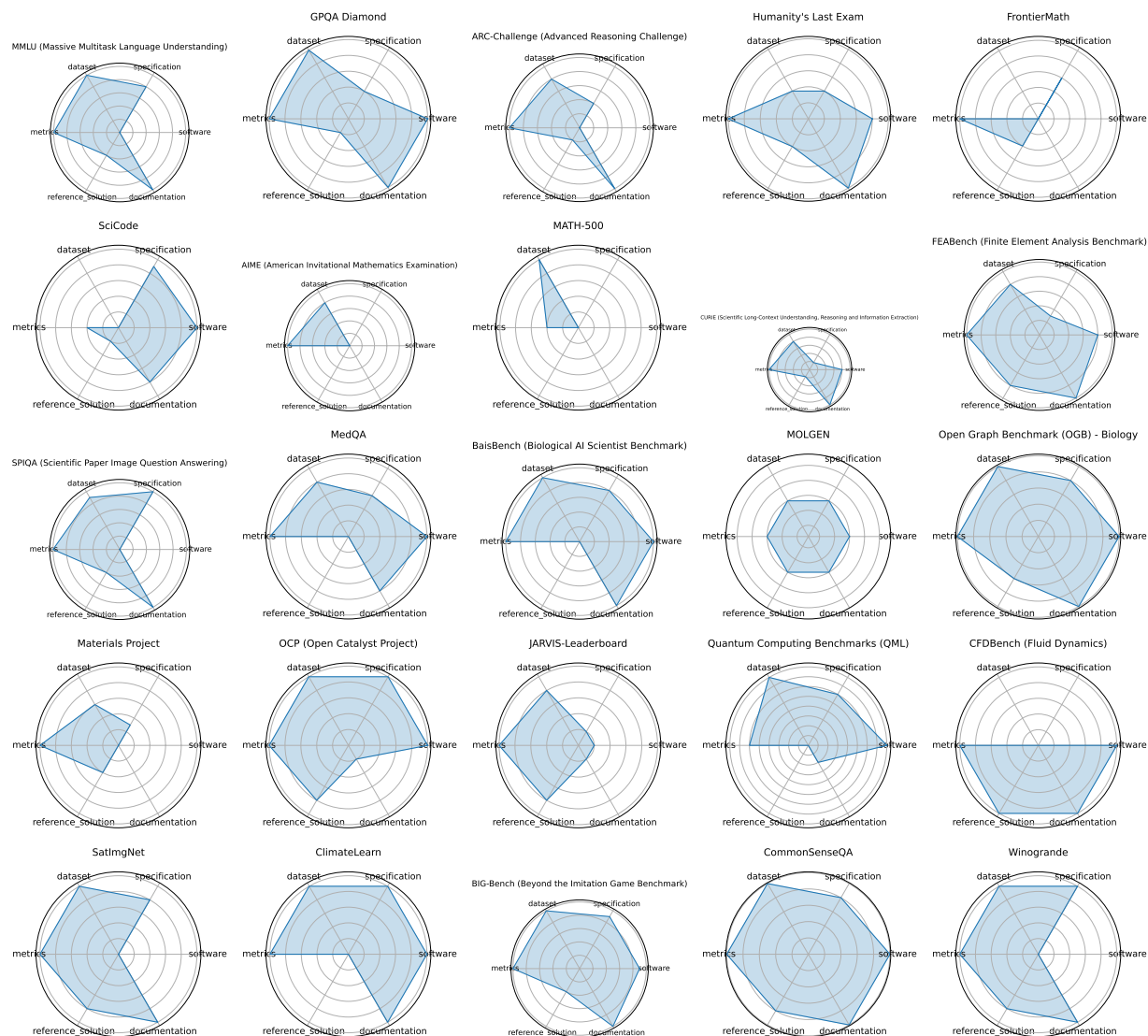


Figure 1: Radar chart overview (page 1)

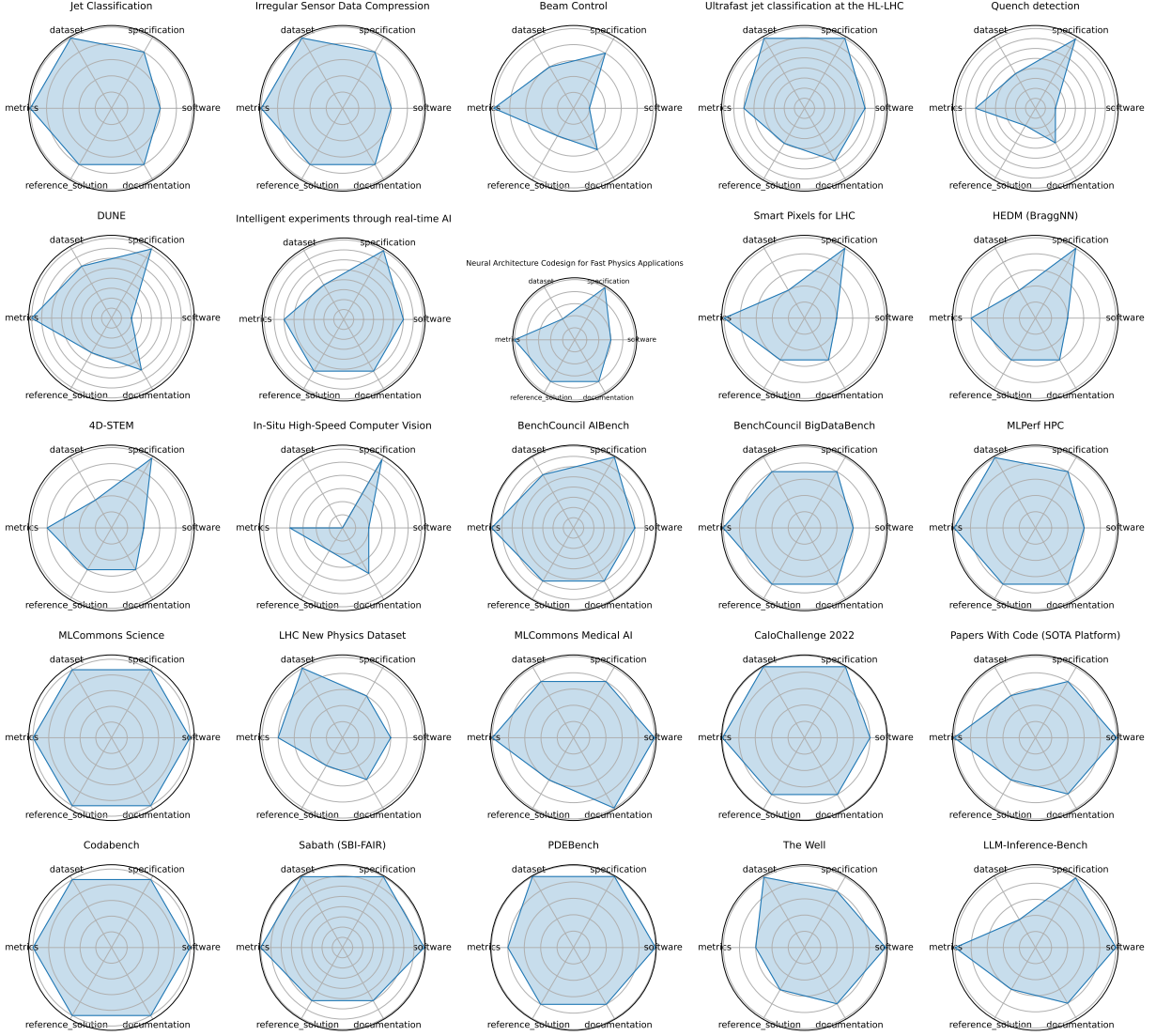


Figure 2: Radar chart overview (page 2)

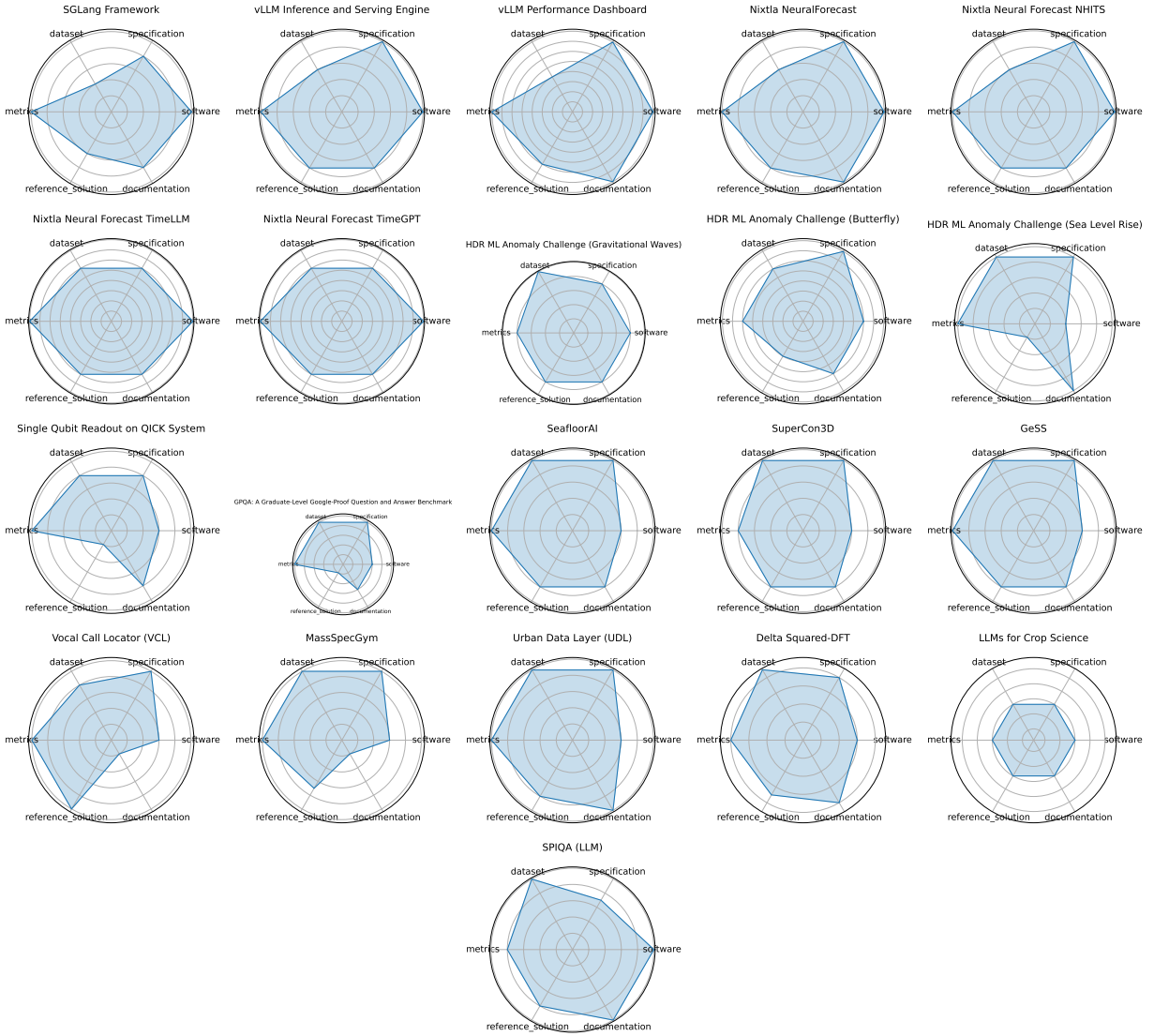


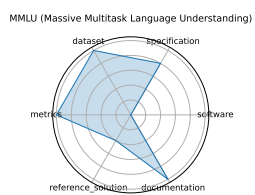
Figure 3: Radar chart overview (page 3)

3 Benchmark Details

4 MMLU (Massive Multitask Language Understanding)

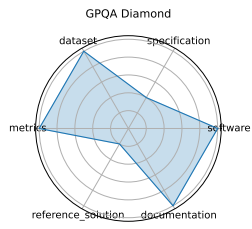
date: 2020-09-07
version: 1
last_updated: 2020-09-07
expired: false
valid: yes
valid_date: 2025-07-28
url: <https://paperswithcode.com/dataset/mmlu>
doi: 10.48550/arXiv.2009.03300
domain: Multidomain
focus: Academic knowledge and reasoning across 57 subjects
keywords: - multitask - multiple-choice - zero-shot - few-shot - knowledge probing
summary: Measuring Massive Multitask Language Understanding (MMLU) is a benchmark of 57 multiple-choice tasks covering elementary mathematics, US history, computer science, law, and more, designed to evaluate a model's breadth and depth of knowledge in zero-shot and few-shot settings.
licensing: MIT License
task_types: - Multiple choice
ai_capability_measured: - General reasoning, subject-matter understanding
metrics: - Accuracy
models: - GPT-4o - Gemini 1.5 Pro - o1 - DeepSeek-R1
ml_motif: - General knowledge
type: Benchmark
ml_task: - Supervised Learning
solutions: 1
notes: Good
contact.name: Dan Hendrycks
contact.email: dan (at) safe.ai
datasets.links.name: Papers with Code datasets
datasets.links.url: <https://github.com/paperswithcode/paperswithcode-data>
results.links.name: Chinchilla
results.links.url: <https://arxiv.org/abs/2203.15556>
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 0
ratings.software.reason: No instructions to download or run data given on the site
ratings.specification.rating: 4
ratings.specification.reason: No system constraints
ratings.dataset.rating: 5
ratings.dataset.reason: Meets all FAIR principles and properly versioned.
ratings.metrics.rating: 5
ratings.metrics.reason: Fully defined, represents a solution's performance.
ratings.reference_solution.rating: 2
ratings.reference_solution.reason: Reference models are available (i.e. GPT-3), but are not trainable or publicly documented
ratings.documentation.rating: 5
ratings.documentation.reason: Well-explained in a provided paper.
id: mmlu_massive_multitask_language_understanding
Citations: [1]

Ratings:



5 GPQA Diamond

date: 2023-11-20
version: 1
last_updated: 2023-11-20
expired: false
valid: yes
valid_date: 2023-11-20
url: <https://arxiv.org/abs/2311.12022>
doi: 10.48550/arXiv.2311.12022
domain: Science
focus: Graduate-level scientific reasoning
keywords: - Google-proof - graduate-level - science QA - chemistry - physics
summary: GPQA is a dataset of 448 challenging, multiple-choice questions in biology, physics, and chemistry, written by domain experts. It is Google-proof - experts score 65% (74% after error correction) while skilled non-experts with web access score only 34%. State-of-the-art LLMs like GPT-4 reach around 39% accuracy.
licensing: unknown
task_types: - Multiple choice - Multi-step QA
ai_capability_measured: - Scientific reasoning, deep knowledge
metrics: - Accuracy
models: - o1 - DeepSeek-R1
ml_motif: - Science and STEM fields
type: Benchmark
ml_task: - Supervised Learning
solutions: 0
notes: Good
contact.name: Julian Michael
contact.email: julianjm@nyu.edu
datasets.links.name: unknown
datasets.links.url: unknown
results.links.name: unknown
results.links.url: unknown
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 5
ratings.software.reason: Python version and requirements specified on Github site
ratings.specification.rating: 2
ratings.specification.reason: No system constraints or I/O specified
ratings.dataset.rating: 5
ratings.dataset.reason: Easily able to access dataset. Comes with predefined splits as mentioned in the paper
ratings.metrics.rating: 5
ratings.metrics.reason: Each question has a correct answer, representing the tested model's performance.
ratings.reference_solution.rating: 1
ratings.reference_solution.reason: Common models such as GPT-3.5 were compared. They are not open and don't provide requirements
ratings.documentation.rating: 5
ratings.documentation.reason: All information is listed in the associated paper
id: gpqa_diamond
Citations: [2]

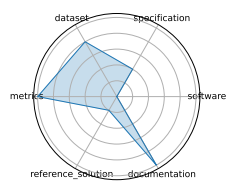


Ratings:

6 ARC-Challenge (Advanced Reasoning Challenge)

date: 2018-03-14
version: 1
last_updated: 2018-03-14
expired: false
valid: yes
valid_date: 2018-03-14
url: <https://allenai.org/data/arc>
doi: NA
domain: Science
focus: Grade-school science with reasoning emphasis
keywords: - grade-school - science QA - challenge set - reasoning
summary: The AI2 Reasoning Challenge (ARC) Challenge set comprises 7,787 natural, grade-school science questions that retrieval-based and word co-occurrence algorithms both fail, requiring advanced reasoning over a 14-million-sentence corpus.
licensing: Apache 2.0 License
task_types: - Multiple choice
ai_capability_measured: - Commonsense and scientific reasoning
metrics: - Accuracy
models: - GPT-4 - Claude
ml_motif: - Elementary science
type: Benchmark
ml_task: - Supervised Learning
solutions: 0
notes: Good
contact.name: unknown
contact.email: unknown
datasets.links.name: Hugging Face
datasets.links.url: https://huggingface.co/datasets/allenai/ai2_arc
results.links.name: unknown
results.links.url: unknown
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 0
ratings.software.reason: No link to code or documentation
ratings.specification.rating: 2
ratings.specification.reason: Task is clear, but no constraints or format is mentioned
ratings.dataset.rating: 4
ratings.dataset.reason: Data accessible, offers instructions on how to download the data via CLI tools. No splits.
ratings.metrics.rating: 5
ratings.metrics.reason: (by default) All questions in the dataset are multiple choice, all have a correct answer
ratings.reference_solution.rating: 1
ratings.reference_solution.reason: There are over 300 models listed, but very few, if any, show performance on the dataset or list constraints
ratings.documentation.rating: 5
ratings.documentation.reason: Explains all necessary information inside a paper
id: arc-challenge_advanced_reasoning_challenge
Citations: [3]

ARC-Challenge (Advanced Reasoning Challenge)

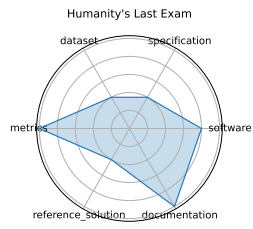


Ratings:

7 Humanity’s Last Exam

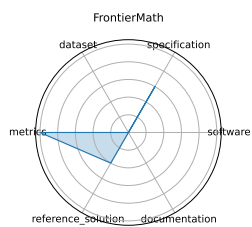
date: 2025-01-24
version: 1
last_updated: 2025-01-24
expired: false
valid: yes
valid_date: 2025-01-24
url: <https://arxiv.org/abs/2501.14249>
doi: 10.48550/arXiv.2501.14249
domain: Multidomain
focus: Broad cross-domain academic reasoning
keywords: - cross-domain - academic exam - multiple-choice - multidisciplinary
summary: Humanity’s Last Exam is a multi-domain, multiple-choice benchmark containing 2,000 questions across diverse academic disciplines, designed to evaluate LLMs’ ability to reason across domains without external resources.
licensing: MIT License
task_types: - Multiple choice
ai_capability_measured: - Cross-domain academic reasoning
metrics: - Accuracy
models: - unknown
ml_motif: - Multi-domain
type: Benchmark
ml_task: - Supervised Learning
solutions: 0
notes: Good
contact.name: HLE team
contact.email: agibenchmark@safe.ai
datasets.links.name: Hugging Face
datasets.links.url: <https://huggingface.co/datasets/cais/hle>
results.links.name: unknown
results.links.url: unknown
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 4
ratings.software.reason: Code for testing models posted on the github. Unknown how to run a custom model.
ratings.specification.rating: 2
ratings.specification.reason: Format of inputs (natural language) and outputs (multiple choice or natural language) specified. No HW constraints specified
ratings.dataset.rating: 2
ratings.dataset.reason: Data accessible through Hugging Face, but requires giving contact information to access
ratings.metrics.rating: 5
ratings.metrics.reason: (by default) All questions in the dataset are multiple choice, all have a correct answer
ratings.reference_solution.rating: 2
ratings.reference_solution.reason: Performance for cutting-edge models listed, but does not specify exact version of the models or how to reproduce the result
ratings.documentation.rating: 5
ratings.documentation.reason: Paper available with necessary information
id: humanitys_last_exam
Citations: [4]

Ratings:



8 FrontierMath

date: 2024-11-07
version: 1
last_updated: 2024-11-07
expired: false
valid: yes
valid_date: 2024-11-07
url: <https://arxiv.org/abs/2411.04872>
doi: 10.48550/arXiv.2411.04872
domain: Mathematics
focus: Challenging advanced mathematical reasoning
keywords: - symbolic reasoning - number theory - algebraic geometry - category theory
summary: FrontierMath is a benchmark of hundreds of expert-vetted mathematics problems spanning number theory, real analysis, algebraic geometry, and category theory, measuring LLMs ability to solve problems requiring deep abstract reasoning.
licensing: unknown
task_types: - Problem solving
ai_capability_measured: - Symbolic and abstract mathematical reasoning
metrics: - Accuracy
models: - unknown
ml_motif: - Math problem solving
type: Benchmark
ml_task: - Supervised Learning
solutions: 0
notes: Good
contact.name: FrontierMath team
contact.email: math_evals@epochai.org
datasets.links.name: unknown
datasets.links.url: unknown
results.links.name: unknown
results.links.url: unknown
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 0
ratings.software.reason: No link to code provided
ratings.specification.rating: 3
ratings.specification.reason: Well-specified process for asking questions and receiving answers. No software or hardware constraints
ratings.dataset.rating: 0
ratings.dataset.reason: Paper and website had no link to any dataset. It may still exist somewhere
ratings.metrics.rating: 5
ratings.metrics.reason: (by default) All questions in the dataset have a correct answer
ratings.reference_solution.rating: 2
ratings.reference_solution.reason: Displays result of leading models on the benchmark, but none are trainable or list constraints
ratings.documentation.rating: 0
ratings.documentation.reason: No specified way to reproduce the reference solution
id: frontiermath
Citations: [5]

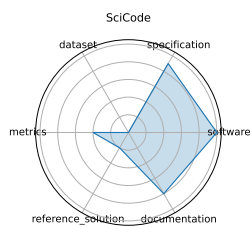


Ratings:

9 SciCode

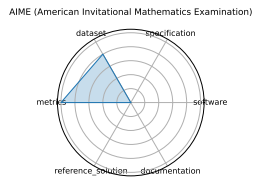
date: 2024-07-18
version: 1
last_updated: 2024-07-18
expired: false
valid: yes
valid_date: 2024-07-18
url: <https://arxiv.org/abs/2407.13168>
doi: 10.48550/arXiv.2407.13168
domain: Scientific Programming
focus: Scientific code generation and problem solving
keywords: - code synthesis - scientific computing - programming benchmark
summary: SciCode is a scientist-curated coding benchmark with 338 subproblems derived from 80 real research tasks across 16 scientific subfields, evaluating models on knowledge recall, reasoning, and code synthesis for scientific computing tasks.
licensing: unknown
task_types: - Coding
ai_capability_measured: - Program synthesis, scientific computing
metrics: - Solve rate (%)
models: - Claude3.5-Sonnet
ml_motif: - Coding
type: Benchmark
ml_task: - Supervised Learning
solutions: unknown
notes: Good
contact.name: Minyang Tian
contact.email: mtian8@illinois.edu
datasets.links.name: unknown
datasets.links.url: unknown
results.links.name: unknown
results.links.url: unknown
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 5
ratings.software.reason: Code to run exists on github repo
ratings.specification.rating: 4.5
ratings.specification.reason: Expected outputs and broad types of inputs stated. Few details on output grading. No HW constraints.
ratings.dataset.rating: 0
ratings.dataset.reason: Paper and website had no link to any dataset. It may still exist somewhere
ratings.metrics.rating: 2
ratings.metrics.reason: Metrics stated, but method of grading is not specified
ratings.reference_solution.rating: 1
ratings.reference_solution.reason: Models presented with scores, but none are open or list constraints
ratings.documentation.rating: 4
ratings.documentation.reason: Paper containing all needed info except for evaluation criteria
id: scicode
Citations: [6]

Ratings:



10 AIME (American Invitational Mathematics Examination)

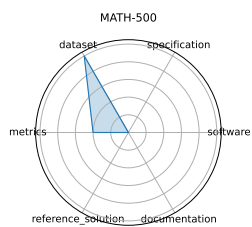
date: 2025-03-13
version: 1
last_updated: 2025-03-13
expired: false
valid: yes
valid_date: 2025-03-13
url: https://artofproblemsolving.com/wiki/index.php/AIME_Problems_and_Solutions
doi: NA
domain: Mathematics
focus: Pre-college advanced problem solving
keywords: - algebra - combinatorics - number theory - geometry
summary: The AIME is a 15-question, 3-hour exam for high-school students featuring challenging short-answer math problems in algebra, number theory, geometry, and combinatorics, assessing depth of problem-solving ability.
licensing: unknown
task_types: - Problem solving
ai_capability_measured: - Mathematical problem-solving and reasoning
metrics: - Accuracy
models: - unknown
ml_motif: - Math problem solving
type: Benchmark
ml_task: - Supervised Learning
solutions: 0
notes: Designed for human test-takers
contact.name: unknown
contact.email: unknown
datasets.links.name: AoPS website
datasets.links.url: https://artofproblemsolving.com/wiki/index.php/AIME_Problems_and_Solutions
results.links.name: unknown
results.links.url: unknown
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 0
ratings.software.reason: No code available
ratings.specification.rating: 0
ratings.specification.reason: Obvious what the problems are, but not specified how to administer them to AI models. No HW constraints
ratings.dataset.rating: 4
ratings.dataset.reason: Easily accessible data with problems and solutions, but no splits
ratings.metrics.rating: 5
ratings.metrics.reason: (by default) Answer is correct or it's not
ratings.reference_solution.rating: 0
ratings.reference_solution.reason: Not given. Human performance stats exist, but no mentions of AI performance
ratings.documentation.rating: 0
ratings.documentation.reason: Not given
id: aime_american_invitational_mathematics_examination
Citations: [7]



Ratings:

11 MATH-500

date: 2025-02-15
version: 1
last_updated: 2025-02-15
expired: false
valid: yes
valid_date: 2025-02-15
url: <https://huggingface.co/datasets/HuggingFaceH4/MATH-500>
doi: unknown
domain: Mathematics
focus: Math reasoning generalization
keywords: - calculus - algebra - number theory - geometry
summary: MATH-500 is a curated subset of 500 problems from the OpenAI MATH dataset, spanning high-school to advanced levels, designed to evaluate LLMs mathematical reasoning and generalization.
licensing: MIT License
task_types: - Problem solving
ai_capability_measured: - Math reasoning and generalization
metrics: - Accuracy
models: - unknown
ml_motif: - Math problem solving
type: Benchmark
ml_task: - Supervised Learning
solutions: 0
notes: Dataset hosted on Hugging Face. Data comes from a subset of OpenAI's dataset
contact.name: unknown
contact.email: unknown
datasets.links.name: Hugging Face
datasets.links.url: <https://huggingface.co/datasets/HuggingFaceH4/MATH-500>
results.links.name: unknown
results.links.url: unknown
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 0
ratings.software.reason: No code provided
ratings.specification.rating: 0
ratings.specification.reason: No method of presentation and evaluation is not stated. No constraints
ratings.dataset.rating: 5
ratings.dataset.reason: Problems and solutions are easily downloaded. Could not find a way to download the data
ratings.metrics.rating: 2
ratings.metrics.reason: Problem spec states that all of the AI reasoning steps are subject to grading, but no specified way to evaluate the steps
ratings.reference_solution.rating: 0
ratings.reference_solution.reason: Not given
ratings.documentation.rating: 0
ratings.documentation.reason: Not given. Implicit instructions to download dataset.
id: math-
Citations: [8]



Ratings:

12 CURIE (Scientific Long-Context Understanding, Reasoning and Information Extraction)

date: 2024-04-02
version: 1
last_updated: 2024-04-02
expired: false
valid: yes
valid_date: 2024-04-02
url: <https://arxiv.org/abs/2503.13517>
doi: 10.48550/arXiv.2503.13517
domain: Multidomain Science
focus: Long-context scientific reasoning
keywords: - long-context - information extraction - multimodal
summary: CURIE is a benchmark of 580 problems across six scientific disciplines-materials science, quantum computing, biology, chemistry, climate science, and astrophysics- designed to evaluate LLMs on long-context understanding, reasoning, and information extraction in realistic scientific workflows.
licensing: Apache 2.0 License
task_types: - Information extraction - Reasoning - Concept tracking - Aggregation - Algebraic manipulation - Multimodal comprehension
ai_capability_measured: - Long-context understanding and scientific reasoning
metrics: - Accuracy
models: - unknown
ml_motif: - Scientific problem solving
type: Benchmark
ml_task: - Supervised Learning
solutions: 0
notes: Good
contact.name: Subhashini Venugopalan
contact.email: vsubhashini@google.com
datasets.links.name: unknown
datasets.links.url: unknown
results.links.name: unknown
results.links.url: unknown
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 4
ratings.software.reason: Code is available, but not well documented
ratings.specification.rating: 1
ratings.specification.reason: Explains types of problems in detail, but does not state exactly how to administer them.
ratings.dataset.rating: 4
ratings.dataset.reason: Dataset is available via Github, but hard to find
ratings.metrics.rating: 5
ratings.metrics.reason: Quantitative metrics such as ROUGE-L and F1 used. Metrics are tailored to the specific problem.
ratings.reference_solution.rating: 1
ratings.reference_solution.reason: Exists, but is not open
ratings.documentation.rating: 5
ratings.documentation.reason: Associated paper explains all criteria
id: curie_scientific_long-context_understanding_reasoning_and_information_extraction
Citations: [9]

CURE (Scientific Long-Context Understanding, Reasoning and Information Extraction)

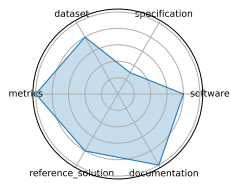


Ratings:

13 FEABench (Finite Element Analysis Benchmark)

date: 2023-01-26
version: 1
last_updated: 2023-01-26
expired: false
valid: no
valid_date: 2023-01-26
url: <https://github.com/google/feabench>
doi: unknown
domain: Computational Engineering
focus: FEA simulation accuracy and performance
keywords: - finite element - simulation - PDE
summary: Does not exist
licensing: unknown
task_types: - Simulation - Performance evaluation
ai_capability_measured: - Numerical simulation accuracy and efficiency
metrics: - Solve time - Error norm
models: - FEniCS - deal.II
ml_motif: - unknown
type: Benchmark
ml_task: - Supervised Learning
solutions: unknown
notes: OK
contact.name: unknown
contact.email: unknown
datasets.links.name: unknown
datasets.links.url: unknown
results.links.name: unknown
results.links.url: unknown
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 4
ratings.software.reason: Code is available, but poorly documented
ratings.specification.rating: 1.5
ratings.specification.reason: Output is defined and task clarity is questionable
ratings.dataset.rating: 4
ratings.dataset.reason: Available, but not split into sets
ratings.metrics.rating: 5
ratings.metrics.reason: Fully defined metrics
ratings.reference_solution.rating: 4
ratings.reference_solution.reason: Three open-source models were used. No system constraints.
ratings.documentation.rating: 5
ratings.documentation.reason: In associated paper
id: feabench_finite_element_analysis_benchmark
Citations: [10]

FEABench (Finite Element Analysis Benchmark)

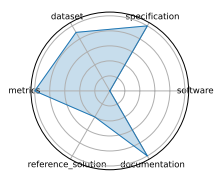


Ratings:

14 SPIQA (Scientific Paper Image Question Answering)

date: 2024-07-12
version: 1
last_updated: 2024-07-12
expired: false
valid: yes
valid_date: 2024-07-12
url: <https://arxiv.org/abs/2407.09413>
doi: 10.48550/arXiv.2407.09413
domain: Computer Science
focus: Multimodal QA on scientific figures
keywords: - multimodal QA - figure understanding - table comprehension - chain-of-thought
summary: SPIQA assesses AI models' ability to interpret and answer questions about figures and tables in scientific papers by integrating visual and textual modalities with chain-of-thought reasoning.
licensing: Apache 2.0 License
task_types: - Question answering - Multimodal QA - Chain-of-Thought evaluation
ai_capability_measured: - Visual-textual reasoning in scientific contexts
metrics: - Accuracy - F1 score
models: - Chain-of-Thought models - Multimodal QA systems
ml_motif: - Scientific paper reading
type: Benchmark
ml_task: - Supervised Learning
solutions: 0
notes: Good
contact.name: Subhashini Venugopalan
contact.email: vsubhashini@google.com
datasets.links.name: Hugging Face
datasets.links.url: <https://huggingface.co/datasets/google/spiqa>
results.links.name: unknown
results.links.url: unknown
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 0
ratings.software.reason: Not provided
ratings.specification.rating: 5
ratings.specification.reason: Task administration clearly defined; prompt instructions explicitly given, no ambiguity in format or scope.
ratings.dataset.rating: 4.5
ratings.dataset.reason: Dataset is available (via paper/appendix), includes train/test/valid split. FAIR-compliant with minor gaps in versioning or access standardization.
ratings.metrics.rating: 5
ratings.metrics.reason: Uses quantitative metrics (Accuracy, F1) aligned with the task
ratings.reference_solution.rating: 2
ratings.reference_solution.reason: Multiple model results (e.g., GPT-4V, Gemini) reported; baselines exist, but full runnable code not confirmed for all.
ratings.documentation.rating: 5
ratings.documentation.reason: All information provided in paper
id: spiqa_scientific_paper_image_question_answering
Citations: [11]

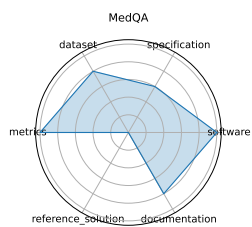
SPIQA (Scientific Paper Image Question Answering)



Ratings:

15 MedQA

date: 2020-09-28
version: 1
last_updated: 2020-09-28
expired: false
valid: yes
valid_date: 2020-09-28
url: <https://arxiv.org/abs/2009.13081>
doi: 10.48550/arXiv.2009.13081
domain: Medical Question Answering
focus: Medical board exam QA
keywords: - USMLE - diagnostic QA - medical knowledge - multilingual
summary: MedQA is a large-scale multiple-choice dataset drawn from professional medical board exams (e.g., USMLE), testing AI systems on diagnostic and medical knowledge questions in English and Chinese.
licensing: Under Association for the Advancement of Artificial Intelligence
task_types: - Multiple choice
ai_capability_measured: - Medical diagnosis and knowledge retrieval
metrics: - Accuracy
models: - Neural reader - Retrieval-based QA systems
ml_motif: - Medical diagnosis
type: Benchmark
ml_task: - Supervised Learning
solutions: 0
notes: Multilingual (English, Simplified and Traditional Chinese)
contact.name: Di Jin
contact.email: jindi15@mit.edu
datasets.links.name: Github
datasets.links.url: <https://github.com/jindi11/MedQA>
results.links.name: unknown
results.links.url: unknown
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 5
ratings.software.reason: All code available on the github
ratings.specification.rating: 3
ratings.specification.reason: Task is clearly defined as multiple-choice QA for medical board exams; input and output formats are explicit; task scope is rigorous and structured. System constraints not specified.
ratings.dataset.rating: 4
ratings.dataset.reason: Dataset is publicly available (GitHub, paper, Hugging Face), well-structured. However, versioning and metadata could be more standardized to fully meet FAIR criteria.
ratings.metrics.rating: 5
ratings.metrics.reason: Uses clear, quantitative metric (accuracy), standard for multiple-choice benchmarks; easily comparable across models.
ratings.reference_solution.rating: 0
ratings.reference_solution.reason: No reference solution mentioned.
ratings.documentation.rating: 4
ratings.documentation.reason: Paper is available. Evaluation criteria are not mentioned.
id: medqa
Citations: [12]

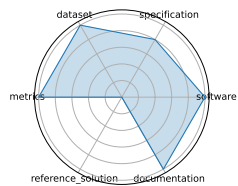


Ratings:

16 BaisBench (Biological AI Scientist Benchmark)

date: 2025-05-13
version: 1
last_updated: 2025-05-13
expired: false
valid: yes
valid_date: 2025-05-13
url: <https://arxiv.org/abs/2505.08341>
doi: 10.48550/arXiv.2505.08341
domain: Computational Biology
focus: Omics-driven AI research tasks
keywords: - single-cell annotation - biological QA - autonomous discovery
summary: BaisBench evaluates AI scientists' ability to perform data-driven biological research by annotating cell types in single-cell datasets and answering MCQs derived from biological study insights, measuring autonomous scientific discovery.
licensing: MIT License
task_types: - Cell type annotation - Multiple choice
ai_capability_measured: - Autonomous biological research capabilities
metrics: - Annotation accuracy - QA accuracy
models: - LLM-based AI scientist agents
ml_motif: - Scientific research
type: Benchmark
ml_task: - Supervised Learning
solutions: 0
notes: Underperforms human experts; aims to advance AI-driven discovery
contact.name: Xuegong Zhang
contact.email: zhangxg@mail.tsinghua.edu.cn
datasets.links.name: Github
datasets.links.url: <https://github.com/EperLuo/BaisBench>
results.links.name: unknown
results.links.url: unknown
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 5
ratings.software.reason: Instructions for environment setup available
ratings.specification.rating: 4
ratings.specification.reason: Task clearly defined-cell type annotation and biological QA; input/output formats are well-described; system constraints are not quantified.
ratings.dataset.rating: 5
ratings.dataset.reason: Uses public scRNA-seq datasets linked in paper appendix; structured and accessible, though versioning and full metadata not formalized per FAIR standards.
ratings.metrics.rating: 5
ratings.metrics.reason: Includes precise and interpretable metrics (annotation and QA accuracy); directly aligned with task outputs and benchmarking goals.
ratings.reference_solution.rating: 0
ratings.reference_solution.reason: Model evaluations and LLM agent results discussed; however, no fully packaged, runnable baseline confirmed yet.
ratings.documentation.rating: 5
ratings.documentation.reason: Dataset and paper accessible; IPYNB files for setup are available on the github repo.
id: baisbench_biological_ai_scientist_benchmark
Citations: [13]

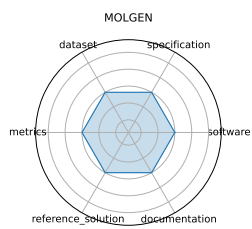
BaisBench (Biological AI Scientist Benchmark)



Ratings:

17 MOLGEN

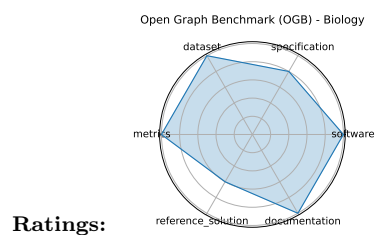
date: 2023-01-26
version: 1
last_updated: 2023-01-26
expired: false
valid: yes
valid_date: 2023-01-26
url: <https://github.com/zjunlp/MolGen>
doi: 10.48550/arXiv.2301.11259
domain: Computational Chemistry
focus: Molecular generation and optimization
keywords: - SELFIES - GAN - property optimization
summary: MolGen is a pre-trained molecular language model that generates chemically valid molecules using SELFIES and reinforcement learning, guided by chemical feedback to optimize properties such as logP, QED, and docking score.
licensing: MIT License
task_types: - Distribution learning - Goal-oriented generation
ai_capability_measured: - Generation of valid and optimized molecular structures
metrics: - Validity% - Novelty% - QED - Docking score
models: - MolGen
ml_motif: - Chemical generation
type: Benchmark
ml_task: - Supervised Learning
solutions: 0
notes: This is a model, not a benchmark
contact.name: unknown
contact.email: unknown
datasets.links.name: unknown
datasets.links.url: unknown
results.links.name: unknown
results.links.url: unknown
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 0
ratings.software.reason: This is a pre-trained model
ratings.specification.rating: 0
ratings.specification.reason: This is a pre-trained model
ratings.dataset.rating: 0
ratings.dataset.reason: This is a pre-trained model
ratings.metrics.rating: 0
ratings.metrics.reason: This is a pre-trained model
ratings.reference_solution.rating: 0
ratings.reference_solution.reason: This is a pre-trained model
ratings.documentation.rating: 0
ratings.documentation.reason: This is a pre-trained model
id: molgen
Citations: [14]



Ratings:

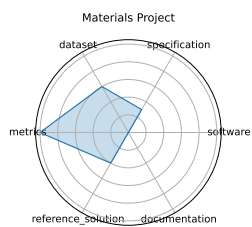
18 Open Graph Benchmark (OGB) - Biology

date: 2020-05-02
version: 1
last_updated: 2020-05-02
expired: false
valid: yes
valid_date: 2020-05-02
url: <https://ogb.stanford.edu/docs/home/>
doi: 10.48550/arXiv.2005.00687
domain: Graph ML
focus: Biological graph property prediction
keywords: - node prediction - link prediction - graph classification
summary: OGB-Biology is a suite of large-scale biological network datasets (protein-protein interaction, drug-target, etc.) with standardized splits and evaluation protocols for node, link, and graph property prediction tasks.
licensing: MIT License
task_types: - Node property prediction - Link property prediction - Graph property prediction
ai_capability_measured: - Scalability and generalization in graph ML for biology
metrics: - Accuracy - ROC-AUC
models: - GCN - GraphSAGE - GAT
ml_motif: - Chemical biology
type: Benchmark
ml_task: - Supervised Learning
solutions: 0
notes: Community-driven updates
contact.name: OGB Team
contact.email: ogb@cs.stanford.edu
datasets.links.name: OGB Webpage
datasets.links.url: https://ogb.stanford.edu/docs/dataset_overview/
results.links.name: unknown
results.links.url: unknown
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 5
ratings.software.reason: All necessary information is provided on the Github
ratings.specification.rating: 4
ratings.specification.reason: Tasks (node/link/graph property prediction) are clearly specified with input/output formats and standardized protocols; constraints (e.g., splits) are well-defined. No constraints.
ratings.dataset.rating: 5
ratings.dataset.reason: Fully FAIR- datasets are versioned, split, and accessible via a standardized API; extensive metadata and documentation are included.
ratings.metrics.rating: 5
ratings.metrics.reason: Reproducible, quantitative metrics (e.g., ROC-AUC, accuracy) that are tightly aligned with the tasks.
ratings.reference_solution.rating: 3
ratings.reference_solution.reason: Multiple baselines implemented and documented (GCN, GAT, GraphSAGE). No constraints.
ratings.documentation.rating: 5
ratings.documentation.reason: All necessary information is included in a paper.
id: open_graph_benchmark_ogb_-_biology
Citations: [15]



19 Materials Project

date: 2011-10-01
version: 1
last_updated: 2011-10-01
expired: false
valid: yes
valid_date: 2011-10-01
url: <https://materialsproject.org/>
doi: unknown
domain: Materials Science
focus: DFT-based property prediction
keywords: - DFT - materials genome - high-throughput
summary: The Materials Project provides an open-access database of computed properties for inorganic materials via high-throughput density functional theory (DFT), accelerating materials discovery.
licensing: <https://next-gen.materialsproject.org/about/terms>
task_types: - Property prediction
ai_capability_measured: - Prediction of inorganic material properties
metrics: - MAE - R^2
models: - Automatminer - Crystal Graph Neural Networks
ml_motif: - Material properties
type: Benchmark
ml_task: - Supervised Learning
solutions: 0
notes: Core component of the Materials Genome Initiative
contact.name: unknown
contact.email: unknown
datasets.links.name: Materials Project Catalysis Explorer
datasets.links.url: <https://next-gen.materialsproject.org/catalysis>
results.links.name: unknown
results.links.url: unknown
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 0
ratings.software.reason: No instructions available
ratings.specification.rating: 1.5
ratings.specification.reason: The platform offers a wide range of material property prediction tasks, but task framing and I/O formats vary by API use and are not always standardized across use cases.
ratings.dataset.rating: 3
ratings.dataset.reason: API key required to access data. No predefined splits.
ratings.metrics.rating: 5
ratings.metrics.reason: Uses numerical metrics like MAE and R^2
ratings.reference_solution.rating: 2
ratings.reference_solution.reason: Numerous models (e.g., Automatminer, CGCNN) trained on the database, but no constraints or documentation listed.
ratings.documentation.rating: 0
ratings.documentation.reason: No explanations or paper provided
id: materials_project
Citations: [16]

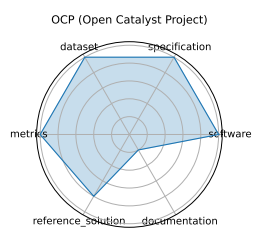


Ratings:

20 OCP (Open Catalyst Project)

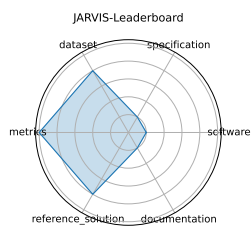
date: 2020-10-20
version: 1
last_updated: 2020-10-20
expired: false
valid: yes
valid_date: 2020-10-20
url: <https://opencatalystproject.org/>
doi: unknown
domain: Chemistry; Materials Science
focus: Catalyst adsorption energy prediction
keywords: - DFT relaxations - adsorption energy - graph neural networks
summary: The Open Catalyst Project (OC20 and OC22) provides DFT-calculated catalyst-adsorbate relaxation datasets, challenging ML models to predict energies and forces for renewable energy applications.
licensing: OCP Terms of Use
task_types: - Energy prediction - Force prediction
ai_capability_measured: - Prediction of adsorption energies and forces
metrics: - MAE (energy) - MAE (force)
models: - CGCNN - SchNet - DimeNet++ - GemNet-OC
ml_motif: - Chemistry
type: Benchmark
ml_task: - Supervised Learning
solutions: 0
notes: Public leaderboards; active community development
contact.name: unknown
contact.email: unknown
datasets.links.name: OCP Dataset
datasets.links.url: <https://fair-chem.github.io/catalysts/datasets/summary>
results.links.name: OCP Pretrained Models
results.links.url: <https://fair-chem.github.io/catalysts/models.html>
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 5
ratings.software.reason: Data provided in Github links
ratings.specification.rating: 5
ratings.specification.reason: Tasks (energy and force prediction) are clearly defined with explicit I/O specifications, constraints, and physical relevance for renewable energy.
ratings.dataset.rating: 5
ratings.dataset.reason: Fully FAIR- OC20, per-adsorbate trajectories, and OC22 are versioned; datasets come with standardized splits, metadata, and are downloadable.
ratings.metrics.rating: 5
ratings.metrics.reason: MAE (energy and force) are standard and reproducible.
ratings.reference_solution.rating: 4
ratings.reference_solution.reason: Multiple baselines (GemNet-OC, DimeNet++, etc.) implemented and evaluated. No hardware listed.
ratings.documentation.rating: 1
ratings.documentation.reason: Paper exists, but content is behind a paywall.
id: ocp_open_catalyst_project
Citations: [17], [18], [19], [20]

Ratings:



21 JARVIS-Leaderboard

date: 2023-06-20
version: 1
last_updated: 2023-06-20
expired: false
valid: yes
valid_date: 2023-06-20
url: <https://arxiv.org/abs/2306.11688>
doi: 10.48550/arXiv.2306.11688
domain: Materials Science; Benchmarking
focus: Comparative evaluation of materials design methods
keywords: - leaderboards - materials methods - simulation
summary: JARVIS-Leaderboard is a community-driven platform benchmarking AI, electronic structure, force-fields, quantum computing, and experimental methods across hundreds of materials science tasks.
licensing: NIST
task_types: - Method benchmarking - Leaderboard ranking
ai_capability_measured: - Performance comparison across diverse materials design methods
metrics: - MAE - RMSE - Accuracy
models: - unknown
ml_motif: - Material science
type: Benchmark
ml_task: - Supervised Learning
solutions: 0
notes: 1281 contributions across 274 benchmarks
contact.name: Kamal Choudhary
contact.email: kamal.choudhary@nist.gov
datasets.links.name: AI model specific benchmarks
datasets.links.url: https://pages.nist.gov/jarvis_leaderboard/AI/
results.links.name: unknown
results.links.url: unknown
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 1
ratings.software.reason: Setup script provided, but no code provided
ratings.specification.rating: 1
ratings.specification.reason: Only dataset format is defined.
ratings.dataset.rating: 4
ratings.dataset.reason: Data is public and adheres to FAIR principles across the NIST-hosted infrastructure; however, metadata completeness varies slightly across benchmarks. No splits.
ratings.metrics.rating: 5
ratings.metrics.reason: Metrics stated for each benchmark.
ratings.reference_solution.rating: 4
ratings.reference_solution.reason: Many baselines across tasks (CGCNN, ALIGNN, M3GNet, etc.); no constraints specified.
ratings.documentation.rating: 1
ratings.documentation.reason: Only the task is specified.
id: jarvis-leaderboard
Citations: [21]

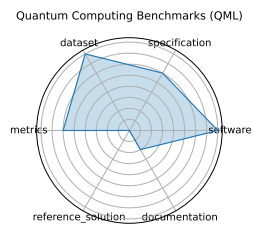


Ratings:

22 Quantum Computing Benchmarks (QML)

date: 2022-02-22
version: 1
last_updated: 2022-02-22
expired: false
valid: yes
valid_date: 2022-02-22
url: <https://github.com/XanaduAI/qml-benchmarks>
doi: 10.48550/arXiv.2307.03901
domain: Quantum Computing
focus: Quantum algorithm performance evaluation
keywords: - quantum circuits - state preparation - error correction
summary: A suite of benchmarks evaluating quantum hardware and algorithms on tasks such as state preparation, circuit optimization, and error correction across multiple platforms.
licensing: Apache-2.0
task_types: - Circuit benchmarking - State classification
ai_capability_measured: - Quantum algorithm performance and fidelity
metrics: - Fidelity - Success probability
models: - IBM Q - IonQ - AQT@LBNL
ml_motif: - Performance Evaluation
type: Benchmark
ml_task: - Supervised Learning
solutions: Varies per benchmark
notes: Hardware-agnostic, application-level metrics. The citation may not be correct.
contact.name: Xanadu AI
contact.email: support@xanadu.ai
datasets.links.name: PennyLane QML Benchmarks Datasets
datasets.links.url: <https://pennylane.ai/datasets/collection/qml-benchmarks>
results.links.name: QML Benchmarks GitHub Repository (Results section)
results.links.url: <https://github.com/XanaduAI/qml-benchmarks#results-and-leaderboards>
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 4
ratings.software.reason: Run instructions exist, but are not easy to follow
ratings.specification.rating: 3
ratings.specification.reason: No system constraints. Task clarity and dataset format are not clearly specified.
ratings.dataset.rating: 4
ratings.dataset.reason: Datasets are accessible, but not split.
ratings.metrics.rating: 3
ratings.metrics.reason: Partially defined, somewhat inferrable metrics. Unknown whether a system's performance is captured.
ratings.reference_solution.rating: 0
ratings.reference_solution.reason: Not provided
ratings.documentation.rating: 1
ratings.documentation.reason: Only the task is defined.
id: quantum_computing_benchmarks_qml
Citations: [22]

Ratings:



23 CFDBench (Fluid Dynamics)

date: 2024-10-01

version: 1

last_updated: 2024-10-01

expired: false

valid: yes

valid_date: 2024-10-01

url: <https://arxiv.org/abs/2310.05963>

doi: 10.48550/arXiv.2310.05963

domain: Fluid Dynamics; Scientific ML

focus: Neural operator surrogate modeling

keywords: - neural operators - CFD - FNO - DeepONet

summary: CFDBench provides large-scale CFD data for four canonical fluid flow problems, assessing neural operators' ability to generalize to unseen PDE parameters and domains.

licensing: CC-BY-4.0

task_types: - Surrogate modeling

ai_capability_measured: - Generalization of neural operators for PDEs

metrics: - L2 error - MAE

models: - FNO - DeepONet - U-Net

ml_motif: - Generalization

type: Benchmark

ml_task: - Supervised Learning

solutions: Numerous, as it's a benchmark for ML models

notes: 302K frames across 739 cases

contact.name: Yining Luo

contact.email: yining.luo@mail.utoronto.ca

datasets.links.name: unknown

datasets.links.url: unknown

results.links.name: unknown

results.links.url: unknown

fair.reproducible: True

fair.benchmark_ready: True

ratings.software.rating: 5

ratings.software.reason: The benchmark provides Python scripts for data loading, preprocessing, and model training/evaluation

ratings.specification.rating: 0

ratings.specification.reason: Not listed

ratings.dataset.rating: 0

ratings.dataset.reason: Not given

ratings.metrics.rating: 5

ratings.metrics.reason: Quantitative metrics (L2 error, MAE, relative error) are clearly defined and align with regression task objectives.

ratings.reference_solution.rating: 5

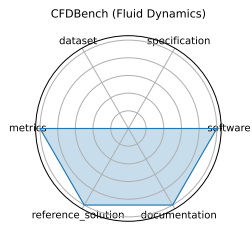
ratings.reference_solution.reason: Baseline models like FNO and DeepONet are implemented, hardware specified.

ratings.documentation.rating: 5

ratings.documentation.reason: Associated paper gives all necessary information.

id: cfdbench_fluid_dynamics

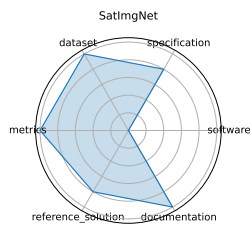
Citations: [23]



Ratings:

24 SatImgNet

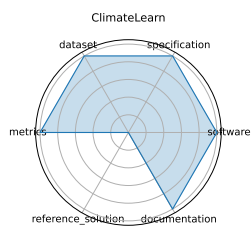
date: 2023-04-23
version: 1
last_updated: 2023-04-23
expired: false
valid: yes
valid_date: 2023-04-23
url: <https://huggingface.co/datasets/saral-ai/satimagnet>
doi: 10.48550/arXiv.2304.11619
domain: Remote Sensing
focus: Satellite imagery classification
keywords: - land-use - zero-shot - multi-task
summary: SATIN (sometimes referred to as SatImgNet) is a multi-task metadataset of 27 satellite imagery classification datasets evaluating zero-shot transfer of vision-language models across diverse remote sensing tasks.
licensing: CC-BY-4.0
task_types: - Image classification
ai_capability_measured: - Zero-shot land-use classification
metrics: - Accuracy
models: - CLIP - BLIP - ALBEF
ml_motif: - Transfer Learning
type: Benchmark
ml_task: - Supervised Learning
solutions: Numerous, evaluated via leaderboard
notes: Public leaderboard available
contact.name: Jonathan Roberts
contact.email: j.roberts@cs.ox.ac.uk
datasets.links.name: SatImgNet on Hugging Face
datasets.links.url: <https://huggingface.co/datasets/saral-ai/satimagnet>
results.links.name: SatImgNet Leaderboard
results.links.url: <https://huggingface.co/spaces/saral-ai/satin-leaderboard>
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 0
ratings.software.reason: No scripts or environment information provided
ratings.specification.rating: 4
ratings.specification.reason: Tasks (image classification across 27 satellite datasets) are clearly defined with multi-task and zero-shot framing; input/output structure is mostly standard but some task-specific nuances require interpretation.
ratings.dataset.rating: 5
ratings.dataset.reason: Hosted on Hugging Face, versioned, FAIR-compliant with rich metadata; covers many well-known remote sensing datasets unified under one metadataset, though documentation depth varies slightly across tasks.
ratings.metrics.rating: 5
ratings.metrics.reason: Accuracy of classification is an appropriate metric
ratings.reference_solution.rating: 4
ratings.reference_solution.reason: Baselines like CLIP, BLIP, ALBEF evaluated in the paper; no constraints specified
ratings.documentation.rating: 5
ratings.documentation.reason: Paper provides all required information
id: satimgnet
Citations: [24]



Ratings:

25 ClimateLearn

date: 2023-07-19
version: 1
last_updated: 2023-07-19
expired: false
valid: yes
valid_date: 2023-07-19
url: <https://arxiv.org/abs/2307.01909>
doi: 10.48550/arXiv.2307.01909
domain: Climate Science; Forecasting
focus: ML for weather and climate modeling
keywords: - medium-range forecasting - ERA5 - data-driven
summary: ClimateLearn provides standardized datasets and evaluation protocols for machine learning models in medium-range weather and climate forecasting using ERA5 reanalysis.
licensing: CC-BY-4.0
task_types: - Forecasting
ai_capability_measured: - Global weather prediction (3-5 days)
metrics: - RMSE - Anomaly correlation
models: - CNN baselines - ResNet variants
ml_motif: - Forecasting - Benchmarking
type: Benchmark
ml_task: - Supervised Learning
solutions: Multiple baseline models provided
notes: Includes physical and ML baselines.
contact.name: Jason Jewik
contact.email: jason.jewik@ucla.edu
datasets.links.name: ClimateLearn GitHub Repository (data loaders and processing)
datasets.links.url: <https://github.com/aditya-grover/climate-learn>
results.links.name: ClimateLearn Paper (results section)
results.links.url: <https://arxiv.org/abs/2307.01909>
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 5
ratings.software.reason: Quickstart notebook makes for easy usage
ratings.specification.rating: 5
ratings.specification.reason: Task framing (medium-range climate forecasting), input/output formats, and evaluation windows are clearly defined; benchmark supports both physical and learned models with detailed constraints.
ratings.dataset.rating: 5
ratings.dataset.reason: Provides standardized access to ERA5 and other reanalysis datasets, with ML-ready splits, meta-data, and Xarray-compatible formats; versioned and fully FAIR-compliant.
ratings.metrics.rating: 5
ratings.metrics.reason: ACC and RMSE are standard, quantitative, and appropriate for climate forecasting; well-integrated into the benchmark, though interpretation across domains may vary.
ratings.reference_solution.rating: 0
ratings.reference_solution.reason: The benchmark is geared for CNN architectures, but no specific model was mentioned.
ratings.documentation.rating: 5
ratings.documentation.reason: Explained in the benchmark's paper.
id: climatelearn
Citations: [25]

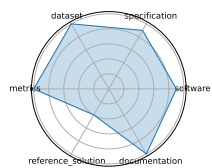


Ratings:

26 BIG-Bench (Beyond the Imitation Game Benchmark)

date: 2022-06-09
version: 1
last_updated: 2022-06-09
expired: false
valid: yes
valid_date: 2022-06-09
url: <https://github.com/google/BIG-bench>
doi: 10.48550/arXiv.2206.04615
domain: NLP; AI Evaluation
focus: Diverse reasoning and generalization tasks
keywords: - few-shot - multi-task - bias analysis
summary: BIG-Bench is a collaborative suite of 204 tasks designed to probe LLMs' reasoning, knowledge, and bias across diverse domains and difficulty levels beyond simple imitation.
licensing: Apache-2.0
task_types: - Few-shot evaluation - Multi-task evaluation
ai_capability_measured: - Reasoning and generalization across diverse tasks
metrics: - Accuracy - Task-specific metrics
models: - GPT-3 - Dense Transformers - Sparse Transformers
ml_motif: - LLM evaluation
type: Benchmark
ml_task: - Supervised Learning
solutions: Multiple, including human baselines
notes: Human baselines included
contact.name: Aarohi Srivastava et al.
contact.email: bigbench@googlegroups.com
datasets.links.name: BIG-Bench GitHub Repository (contains tasks and data)
datasets.links.url: https://github.com/google/BIG-bench/tree/main/bigbench/benchmark_tasks
results.links.name: BIG-Bench GitHub Repository (results in papers and code)
results.links.url: <https://github.com/google/BIG-bench>
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 4.5
ratings.software.reason: Quick start notebook provided, but instructions on how to run it are lacking.
ratings.specification.rating: 4.5
ratings.specification.reason: Tasks are diverse and clearly described; input/output formats are usually defined but vary widely, and system constraints are not standardized.
ratings.dataset.rating: 5
ratings.dataset.reason: Public, versioned, and well-documented; FAIR overall
ratings.metrics.rating: 5
ratings.metrics.reason: Many tasks use standard quantitative metrics (accuracy, BLEU, F1). Others involve subjective ratings (e.g., Likert), which reduces cross-task comparability.
ratings.reference_solution.rating: 2
ratings.reference_solution.reason: Human baselines and LLM performance results are included; however, runnable reference solutions are limited and setup is not fully turnkey.
ratings.documentation.rating: 5
ratings.documentation.reason: Explained in the associated paper.
id: big-bench_beyond_the_imitation_game_benchmark
Citations: [26]

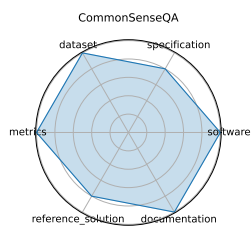
BIG-Bench (Beyond the Imitation Game Benchmark)



Ratings:

27 CommonSenseQA

date: 2019-11-20
version: 1
last_updated: 2019-11-20
expired: false
valid: yes
valid_date: 2019-11-20
url: <https://paperswithcode.com/paper/commonsenseqa-a-question-answering-challenge>
doi: 10.48550/arXiv.1811.00937
domain: NLP; Commonsense
focus: Commonsense question answering
keywords: - ConceptNet - multiple-choice - adversarial
summary: CommonsenseQA is a challenging multiple-choice QA dataset built from ConceptNet, requiring models to apply commonsense knowledge to select the correct answer among five choices.
licensing: MIT
task_types: - Multiple choice
ai_capability_measured: - Commonsense reasoning and knowledge integration
metrics: - Accuracy
models: - BERT-large - RoBERTa - GPT-3
ml_motif: - Commonsense question answering
type: Benchmark
ml_task: - Supervised Learning
solutions: 2
notes: Baseline 56%, human 89%
contact.name: Alon Talmor, Jonathan Herzig, Nicholas Lourie, Jonathan Berant
contact.email: Unknown
datasets.links.name: CommonsenseQA Dataset (Hugging Face)
datasets.links.url: https://huggingface.co/datasets/commonsense_qa
results.links.name: Papers With Code Leaderboard for CommonsenseQA
results.links.url: <https://paperswithcode.com/dataset/commonsenseqa>
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 5
ratings.software.reason: All code given on Github site
ratings.specification.rating: 4
ratings.specification.reason: Task and format (multiple-choice QA with 5 options) are clearly defined; grounded in ConceptNet with consistent structure, though no hardware/system constraints are specified.
ratings.dataset.rating: 5
ratings.dataset.reason: Public, versioned, and FAIR-compliant; includes metadata, splits, and licensing; well-integrated with HuggingFace and other ML libraries.
ratings.metrics.rating: 5
ratings.metrics.reason: Accuracy is a simple, reproducible metric aligned with task goals; no ambiguity in evaluation.
ratings.reference_solution.rating: 4
ratings.reference_solution.reason: Several baseline models (e.g., BERT, RoBERTa) are reported with scores; implementations exist in public repos, but not run with hardware constraints
ratings.documentation.rating: 5
ratings.documentation.reason: Given in paper.
id: commonsenseqa
Citations: [27]

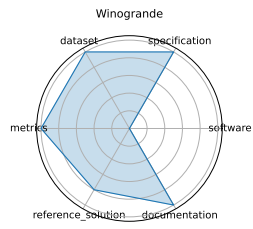


Ratings:

28 Winogrande

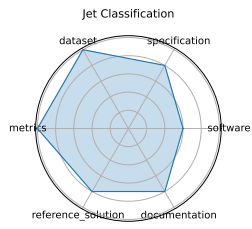
date: 2019-07-24
version: 1
last_updated: 2019-07-24
expired: false
valid: yes
valid_date: 2019-07-24
url: <https://leaderboard.allenai.org/winogrande/submissions/public>
doi: 10.48550/arXiv.1907.10641
domain: NLP; Commonsense
focus: Winograd Schema-style pronoun resolution
keywords: - adversarial - pronoun resolution
summary: WinoGrande is a large-scale adversarial dataset of 44,000 Winograd Schema-style questions with reduced bias using AFLite, serving as both a benchmark and transfer learning resource.
licensing: CC-BY
task_types: - Pronoun resolution
ai_capability_measured: - Robust commonsense reasoning
metrics: - Accuracy - AUC
models: - RoBERTa - BERT - GPT-2
ml_motif: - Commonsense reasoning
type: Benchmark
ml_task: - Supervised Learning
solutions: 2
notes: Human ~94%
contact.name: Keisuke Sakaguchi
contact.email: keisukes@allenai.org
datasets.links.name: Hugging Face / AllenAI
datasets.links.url: <https://huggingface.co/datasets/allenai/winogrande>
results.links.name: Papers With Code leaderboard
results.links.url: <https://paperswithcode.com/dataset/winogrande>
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 0
ratings.software.reason: No template code provided
ratings.specification.rating: 5
ratings.specification.reason: Task (pronoun/coreference resolution) is clearly defined in Winograd Schema style, with consistent input/output format; no system constraints included.
ratings.dataset.rating: 5
ratings.dataset.reason: Public, versioned, and FAIR-compliant with AFLite-generated splits to reduce annotation artifacts; hosted by AllenAI with good metadata.
ratings.metrics.rating: 5
ratings.metrics.reason: Accuracy and AUC are quantitative and well-aligned with disambiguation goals; standardized across evaluations.
ratings.reference_solution.rating: 4
ratings.reference_solution.reason: Baseline results available, requiring users to submit their methods along with their submissions. Constraints are not required in submissions.
ratings.documentation.rating: 5
ratings.documentation.reason: Dataset page and paper provide sufficient detail
id: winogrande
Citations: [28]

Ratings:



29 Jet Classification

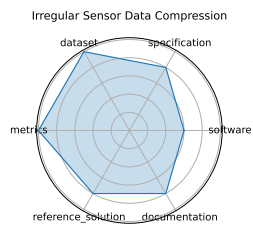
date: 2024-05-01
version: v0.2.0
last_updated: 2024-05
expired: unknown
valid: yes
valid_date: 2024-05-01
url: <https://github.com/fastmachinelearning/fastml-science/tree/main/jet-classify>
doi: 10.48550/arXiv.2207.07958
domain: Particle Physics
focus: Real-time classification of particle jets using HL-LHC simulation features
keywords: - classification - real-time ML - jet tagging - QKeras
summary: This benchmark evaluates ML models for real-time classification of particle jets using high-level features derived from simulated LHC data. It includes both full-precision and quantized models optimized for FPGA deployment.
licensing: Apache License 2.0
task_types: - Classification
ai_capability_measured: - Real-time inference - model compression performance
metrics: - Accuracy - AUC
models: - Keras DNN - QKeras quantized DNN
ml_motif: - Real-time
type: Benchmark
ml_task: - Supervised Learning
solutions: Solution details are described in the referenced paper or repository.
notes: Includes both float and quantized models using QKeras
contact.name: Jules Muhizi
contact.email: unknown
datasets.links.name: JetClass
datasets.links.url: <https://zenodo.org/record/6619768>
results.links.name: ChatGPT LLM
results.links.url: https://docs.google.com/document/d/1runrcij-eoH3_lgGZ8wm2z1YbL1Qf5cSNbVbHyWFDs4
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 3
ratings.software.reason: Not containerized; Setup automation/documentation could be improved
ratings.specification.rating: 4
ratings.specification.reason: System constraints missing
ratings.dataset.rating: 5
ratings.dataset.reason: None
ratings.metrics.rating: 5
ratings.metrics.reason: None
ratings.reference_solution.rating: 4
ratings.reference_solution.reason: HW/SW requirements missing; Reference not bundled as official starter kit
ratings.documentation.rating: 4
ratings.documentation.reason: Full reproducibility requires manual setup
id: jet_classification
Citations: [29]



Ratings:

30 Irregular Sensor Data Compression

date: 2024-05-01
version: v0.2.0
last_updated: 2024-05
expired: unknown
valid: yes
valid_date: 2024-05-01
url: <https://github.com/fastmachinelearning/fastml-science/tree/main/sensor-data-compression>
doi: 10.48550/arXiv.2207.07958
domain: Particle Physics
focus: Real-time compression of sparse sensor data with autoencoders
keywords: - compression - autoencoder - sparse data - irregular sampling
summary: This benchmark addresses lossy compression of irregularly sampled sensor data from particle detectors using real-time autoencoder architectures, targeting latency-critical applications in physics experiments.
licensing: Apache License 2.0
task_types: - Compression
ai_capability_measured: - Reconstruction quality - compression efficiency
metrics: - MSE - Compression ratio
models: - Autoencoder - Quantized autoencoder
ml_motif: - Real-time, Image/CV
type: Benchmark
ml_task: - Unsupervised Learning
solutions: Solution details are described in the referenced paper or repository.
notes: Based on synthetic but realistic physics sensor data
contact.name: Ben Hawks, Nhan Tran
contact.email: unknown
datasets.links.name: Custom synthetic irregular sensor dataset
datasets.links.url: <https://github.com/fastmachinelearning/fastml-science/tree/main/sensor-data-compression>
results.links.name: ChatGPT LLM
fair.reproducible: True
fair.benchmark_ready: True
ratings.software.rating: 3
ratings.software.reason: Not containerized; Full automation and documentation could be improved
ratings.specification.rating: 4
ratings.specification.reason: Exact latency or resource constraints not numerically specified
ratings.dataset.rating: 5
ratings.dataset.reason: All criteria met
ratings.metrics.rating: 5
ratings.metrics.reason: All criteria met
ratings.reference_solution.rating: 4
ratings.reference_solution.reason: Not fully documented or automated for reproducibility
ratings.documentation.rating: 4
ratings.documentation.reason: Setup for deployment (e.g., FPGA pipeline) requires familiarity with tooling
id: irregular_sensor_data_compression
Citations: [30]



Ratings:

31 Beam Control

date: 2024-05-01

version: v0.2.0

last_updated: 2024-05

expired: unknown

valid: yes

valid_date: 2024-05-01

url: <https://github.com/fastmachinelearning/fastml-science/tree/main/beam-control>

doi: 10.48550/arXiv.2207.07958

domain: Accelerators and Magnets

focus: Reinforcement learning control of accelerator beam position

keywords: - RL - beam stabilization - control systems - simulation

summary: Beam Control explores real-time reinforcement learning strategies for maintaining stable beam trajectories in particle accelerators. The benchmark is based on the BOOSTR environment for accelerator simulation.

licensing: Apache License 2.0

task_types: - Control

ai_capability_measured: - Policy performance in simulated accelerator control

metrics: - Stability - Control loss

models: - DDPG - PPO (planned)

ml_motif: - Real-time, RL

type: Benchmark

ml_task: - Reinforcement Learning

solutions: Solution details are described in the referenced paper or repository.

notes: Environment defined, baseline RL implementation is in progress

contact.name: Ben Hawks, Nhan Tran

contact.email: unknown

results.links.name: ChatGPT LLM

fair.reproducible: in progress

fair.benchmark_ready: in progress

ratings.software.rating: 1

ratings.software.reason: Code not documented; Incomplete setup and not containerized

ratings.specification.rating: 4

ratings.specification.reason: Latency/resource constraints not fully quantified

ratings.dataset.rating: 3

ratings.dataset.reason: Not findable (no DOI/indexing); Not interoperable (format/schema unspecified)

ratings.metrics.rating: 5

ratings.metrics.reason: All criteria met

ratings.reference_solution.rating: 2

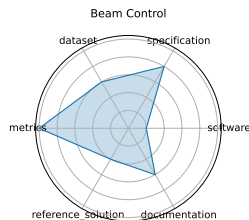
ratings.reference_solution.reason: HW/SW requirements missing; Metrics not evaluated with reference; Baseline not trainable/open

ratings.documentation.rating: 3

ratings.documentation.reason: Setup instructions and pretrained model details are missing

id: beam_control

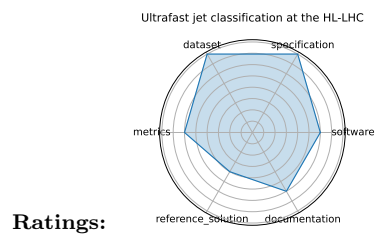
Citations: [31], [32]



Ratings:

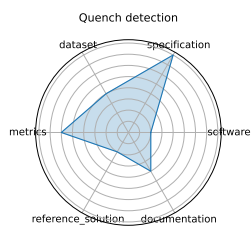
32 Ultrafast jet classification at the HL-LHC

date: 2024-07-08
version: v1.0
last_updated: 2024-07
expired: unknown
valid: yes
valid_date: 2024-07-08
url: <https://arxiv.org/pdf/2402.01876>
doi: 10.48550/arXiv.2402.01876
domain: Particle Physics
focus: FPGA-optimized real-time jet origin classification at the HL-LHC
keywords: - jet classification - FPGA - quantization-aware training - Deep Sets - Interaction Networks
summary: Demonstrates three ML models (MLP, Deep Sets, Interaction Networks) optimized for FPGA deployment with O(100 ns) inference using quantized models and hls4ml, targeting real-time jet tagging in the L1 trigger environment at the high-luminosity LHC. Data is available on Zenodo DOI:10.5281/zenodo.3602260.
licensing: CC-BY
task_types: - Classification
ai_capability_measured: - Real-time inference under FPGA constraints
metrics: - Accuracy - Latency - Resource utilization
models: - MLP - Deep Sets - Interaction Network
ml_motif: - Real-time
type: Model
ml_task: - Supervised Learning
solutions: Solution details are described in the referenced paper or repository.
notes: Uses quantization-aware training; hardware synthesis evaluated via hls4ml
contact.name: Patrick Odagiu
contact.email: podagiu@ethz.ch
datasets.links.name: Zenodo dataset
datasets.links.url: <https://zenodo.org/records/3602260>
results.links.name: ChatGPT LLM
results.links.url: https://docs.google.com/document/d/1gDf1CIYtfmfZ9urv1jCRZMYz_3WwEETkugUC65OZBdw
fair.reproducible: True
fair.benchmark_ready: False
ratings.software.rating: 3
ratings.software.reason: Not containerized; Setup and automation incomplete
ratings.specification.rating: 4
ratings.specification.reason: Hardware constraints are referenced but not fully detailed or standardized
ratings.dataset.rating: 4
ratings.dataset.reason: FAIR metadata limited; no clear mention of dataset format or splits
ratings.metrics.rating: 3
ratings.metrics.reason: Metrics exist (accuracy, latency, utilization), but formal definitions and evaluation guidance are limited
ratings.reference_solution.rating: 2
ratings.reference_solution.reason: Reference implementations not fully reproducible; no evaluation pipeline or training setup provided
ratings.documentation.rating: 3
ratings.documentation.reason: No linked GitHub repo or setup instructions; paper provides partial guidance only
id: ultrafast_jet_classification_at_the_hl-lhc
Citations: [33]



33 Quench detection

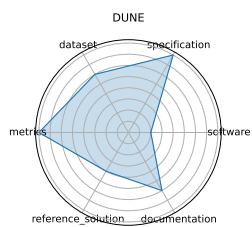
date: 2024-10-15
version: v1.0
last_updated: 2024-10
expired: no
valid: yes
valid_date: 2024-10-15
url: https://indico.cern.ch/event/1387540/contributions/6153618/attachments/2948441/5182077/fast_ml_magnets_2024_final.pdf
doi: NA
domain: Accelerators and Magnets
focus: Real-time detection of superconducting magnet quenches using ML
keywords: - quench detection - autoencoder - anomaly detection - real-time
summary: Exploration of real-time quench detection using unsupervised and RL approaches, combining multi-modal sensor data (BPM, power supply, acoustic), operating on kHz-MHz streams with anomaly detection and frequency-domain features.
licensing: Via Fermilab
task_types: - Anomaly detection - Quench localization
ai_capability_measured: - Real-time anomaly detection with multi-modal sensors
metrics: - ROC-AUC - Detection latency
models: - Autoencoder - RL agents (in development)
ml_motif: - Real-time, RL
type: Benchmark
ml_task: - Reinforcement + Unsupervised Learning
solutions: 0
notes: Precursor detection in progress; multi-modal and dynamic weighting methods
contact.name: Maira Khan
contact.email: unknown
datasets.links.name: BPM and power supply data from BNL
results.links.name: ChatGPT LLM
fair.reproducible: in progress
fair.benchmark_ready: False
ratings.software.rating: 1
ratings.software.reason: Code not provided; no evidence of documentation or containerization
ratings.specification.rating: 4
ratings.specification.reason: Real-time detection task is clearly described, but exact constraints, inputs/outputs, and evaluation protocol are only partially specified
ratings.dataset.rating: 2
ratings.dataset.reason: Dataset URL is missing; FAIR principles largely unmet
ratings.metrics.rating: 3
ratings.metrics.reason: ROC-AUC and latency are mentioned, but metric definitions and formal evaluation setup are missing
ratings.reference_solution.rating: 1
ratings.reference_solution.reason: No baseline or reproducible model implementation available
ratings.documentation.rating: 2
ratings.documentation.reason: Only a conference slide deck is available; lacks detailed instructions or repository for reproduction
id: quench_detection
Citations: [34]



Ratings:

34 DUNE

date: 2024-10-15
version: v1.0
last_updated: 2024-10
expired: unknown
valid: yes
valid_date: 2024-10-15
url: https://indico.fnal.gov/event/66520/contributions/301423/attachments/182439/250508/fast_ml_dunedaq_sonic_10_15_24.pdf
doi: 10.48550/arXiv.2103.13910
domain: Particle Physics
focus: Real-time ML for DUNE DAQ time-series data
keywords: - DUNE - time-series - real-time - trigger
summary: Applying real-time ML methods to time-series data from DUNE detectors, exploring trigger-level anomaly detection and event selection with low latency constraints.
licensing: Via Fermilab
task_types: - Trigger selection - Time-series anomaly detection
ai_capability_measured: - Low-latency event detection
metrics: - Detection efficiency - Latency
models: - CNN - LSTM (planned)
ml_motif: - Real-time, Time-series
type: Benchmark (in progress)
ml_task: - Supervised Learning
solutions: Solution details are described in the referenced paper or repository.
notes: Prototype models demonstrated on SONIC platform
contact.name: Andrew J. Morgan
contact.email: unknown
datasets.links.name: DUNE SONIC data
results.links.name: ChatGPT LLM
fair.reproducible: in progress
fair.benchmark_ready: False
ratings.software.rating: 1
ratings.software.reason: Code not available; no containerization or setup provided
ratings.specification.rating: 4
ratings.specification.reason: Constraints like latency thresholds are described qualitatively but not numerically defined
ratings.dataset.rating: 3
ratings.dataset.reason: Dataset lacks a public URL; FAIR metadata and versioning are missing
ratings.metrics.rating: 4
ratings.metrics.reason: Metrics are relevant but no benchmark baseline or detailed evaluation guidance is provided
ratings.reference_solution.rating: 2
ratings.reference_solution.reason: Autoencoder prototype exists but is not reproducible; RL model still in development
ratings.documentation.rating: 3
ratings.documentation.reason: Documentation exists only in slides/GDocs; no implementation guide or structured release
id: dune
Citations: [35]



Ratings:

35 Intelligent experiments through real-time AI

date: 2025-01-08

version: v1.0

last_updated: 2025-01

expired: unknown

valid: yes

valid_date: 2025-01-08

url: <https://arxiv.org/pdf/2501.04845>

doi: 10.48550/arXiv.2501.04845

domain: Instrumentation and Detectors; Nuclear Physics; Particle Physics

focus: Real-time FPGA-based triggering and detector control for sPHENIX and future EIC

keywords: - FPGA - Graph Neural Network - hls4ml - real-time inference - detector control

summary: Research and Development demonstrator for real-time processing of high-rate tracking data from the sPHENIX detector (RHIC) and future EIC systems. Uses GNNs with hls4ml for FPGA-based trigger generation to identify rare events (heavy flavor, DIS electrons) within 10 micros latency. Demonstrated improved accuracy and latency on Alveo/FELIX platforms.

licensing: CC BY-NC-ND 4.0

task_types: - Trigger classification - Detector control - Real-time inference

ai_capability_measured: - Low-latency GNN inference on FPGA

metrics: - Accuracy (charm and beauty detection) - Latency (micros) - Resource utilization (LUT/FF/BRAM/DSP)

models: - Bipartite Graph Network with Set Transformers (BGN-ST) - GarNet (edge-classifier)

ml_motif: - Real-time

type: Model

ml_task: - Supervised Learning

solutions: Solution details are described in the referenced paper or repository.

notes: Achieved ~97.4% accuracy for beauty decay triggers; sub-10 micros latency on Alveo U280; hit-based FPGA design via hls4ml and FlowGNN.

contact.name: Jakub Kvapil

contact.email: Jakub.Kvapil@lanl.gov

datasets.links.name: Internal simulated tracking data (sPHENIX and EIC DIS-electron tagger)

results.links.name: ChatGPT LLM

fair.reproducible: True

fair.benchmark_ready: False

ratings.software.rating: 3

ratings.software.reason: No containerized or open-source setup provided

ratings.specification.rating: 4

ratings.specification.reason: Architectural/system specifications are incomplete

ratings.dataset.rating: 2

ratings.dataset.reason: Dataset is internal and not publicly available or FAIR-compliant

ratings.metrics.rating: 3

ratings.metrics.reason: Metrics relevant but not supported by evaluation scripts or baselines

ratings.reference_solution.rating: 3

ratings.reference_solution.reason: No public or reproducible implementation released

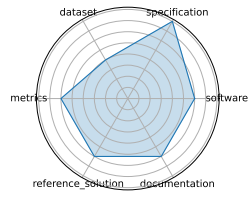
ratings.documentation.rating: 3

ratings.documentation.reason: No public GitHub or complete pipeline documentation

id: intelligent_experiments_through_real-time_ai

Citations: [36]

Intelligent experiments through real-time AI



Ratings:

36 Neural Architecture Codesign for Fast Physics Applications

date: 2025-01-09

version: v1.0

last_updated: 2025-01

expired: unknown

valid: yes

valid_date: 2025-01-09

url: <https://arxiv.org/abs/2501.05515>

doi: 10.48550/arXiv.2501.05515

domain: Physics; Materials Science; Particle Physics

focus: Automated neural architecture search and hardware-efficient model codesign for fast physics applications

keywords: - neural architecture search - FPGA deployment - quantization - pruning - hls4ml

summary: Introduces a two-stage neural architecture codesign (NAC) pipeline combining global and local search, quantization-aware training, and pruning to design efficient models for fast Bragg peak finding and jet classification, synthesized for FPGA deployment with hls4ml. Achieves >30x reduction in BOPs and sub-100 ns inference latency on FPGA.

licensing: Via Fermilab

task_types: - Classification - Peak finding

ai_capability_measured: - Hardware-aware model optimization; low-latency inference

metrics: - Accuracy - Latency - Resource utilization

models: - NAC-based BraggNN - NAC-optimized Deep Sets (jet)

ml_motif: - Real-time, Image/CV

type: Framework

ml_task: - Supervised Learning

solutions: Solution details are described in the referenced paper or repository.

notes: Demonstrated two case studies (materials science, HEP); pipeline and code open-sourced.

contact.name: Jason Weitz (UCSD), Nhan Tran (FNAL)

contact.email: unknown

results.links.name: ChatGPT LLM

fair.reproducible: Yes (nac-opt, hls4ml)

fair.benchmark_ready: False

ratings.software.rating: 3

ratings.software.reason: Toolchain (hls4ml, nac-opt) described but not yet containerized or fully packaged

ratings.specification.rating: 5

ratings.specification.reason: Fully specified task with constraints and target deployment; includes hardware context

ratings.dataset.rating: 2

ratings.dataset.reason: Simulated datasets referenced but not publicly available or FAIR-compliant

ratings.metrics.rating: 5

ratings.metrics.reason: Clear, quantitative metrics aligned with task goals and hardware evaluation

ratings.reference_solution.rating: 4

ratings.reference_solution.reason: Models tested on hardware with source code references; full training pipeline not yet released

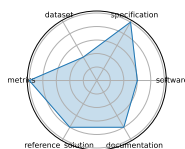
ratings.documentation.rating: 4

ratings.documentation.reason: Detailed paper and tools described; open repo planned but not yet complete

id: neural_architecture_codesign_for_fast_physics_applications

Citations: [37]

Neural Architecture Codesign for Fast Physics Applications



Ratings:

37 Smart Pixels for LHC

date: 2024-06-24

version: v1.0

last_updated: 2024-06

expired: unknown

valid: yes

valid_date: 2024-06-24

url: <https://arxiv.org/abs/2406.14860>

doi: 10.48550/arXiv.2406.14860

domain: Particle Physics; Instrumentation and Detectors

focus: On-sensor, in-pixel ML filtering for high-rate LHC pixel detectors

keywords: - smart pixel - on-sensor inference - data reduction - trigger

summary: Presents a 256x256-pixel ROIC in 28 nm CMOS with embedded 2-layer NN for cluster filtering at 25 ns, achieving 54-75% data reduction while maintaining noise and latency constraints. Prototype consumes ~300 microW/pixel and operates in combinatorial digital logic.

licensing: Via Fermilab

task_types: - Image Classification - Data filtering

ai_capability_measured: - On-chip - low-power inference; data reduction

metrics: - Data rejection rate - Power per pixel

models: - 2-layer pixel NN

ml_motif: - Real-time, Image/CV

type: Benchmark

ml_task: - Image Classification

solutions: Solution details are described in the referenced paper or repository.

notes: Prototype in CMOS 28 nm; proof-of-concept for Phase III pixel upgrades.

contact.name: Lindsey Gray; Jennet Dickinson

contact.email: unknown

results.links.name: ChatGPT LLM

fair.reproducible: True

fair.benchmark_ready: Yes (Zenodo:7331128)

ratings.software.rating: 2

ratings.software.reason: No packaged code or setup scripts available; replication depends on hardware description and paper

ratings.specification.rating: 5

ratings.specification.reason: None

ratings.dataset.rating: 2

ratings.dataset.reason: No dataset links; not publicly hosted or FAIR-compliant

ratings.metrics.rating: 5

ratings.metrics.reason: None

ratings.reference_solution.rating: 3

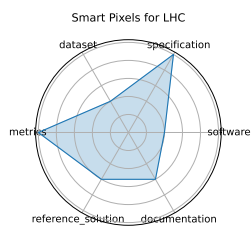
ratings.reference_solution.reason: In-pixel 2-layer NN described and evaluated, but reproducibility and source files are not released

ratings.documentation.rating: 3

ratings.documentation.reason: Paper contains detailed descriptions, but no repo or external guide for reproducing results

id: smart_pixels_for_lhc

Citations: [38]



Ratings:

38 HEDM (BraggNN)

date: 2023-10-03

version: v1.0

last_updated: 2023-10

expired: unknown

valid: yes

valid_date: 2023-10-03

url: <https://arxiv.org/abs/2008.08198>

doi: 10.48550/arXiv.2008.08198

domain: Material Science

focus: Fast Bragg peak analysis using deep learning in diffraction microscopy

keywords: - BraggNN - diffraction - peak finding - HEDM

summary: Uses BraggNN, a deep neural network, for rapid Bragg peak localization in high-energy diffraction microscopy, achieving about 13x speedup compared to Voigt-based methods while maintaining sub-pixel accuracy.

licensing: DOE Public Access Plan

task_types: - Peak detection

ai_capability_measured: - High-throughput peak localization

metrics: - Localization accuracy - Inference time

models: - BraggNN

ml_motif: - Real-time, Image/CV

type: Framework

ml_task: - Peak finding

solutions: Solution details are described in the referenced paper or repository.

notes: Enables real-time HEDM workflows; basis for NAC case study.

contact.name: Jason Weitz (UCSD)

contact.email: unknown

results.links.name: ChatGPT LLM

fair.reproducible: True

fair.benchmark_ready: True

ratings.software.rating: 2

ratings.software.reason: No standalone code repository or setup instructions provided

ratings.specification.rating: 5

ratings.specification.reason: None

ratings.dataset.rating: 2

ratings.dataset.reason: No dataset links or FAIR metadata; unclear public access

ratings.metrics.rating: 4

ratings.metrics.reason: Only localization accuracy and inference time mentioned; not formally benchmarked with scripts

ratings.reference_solution.rating: 3

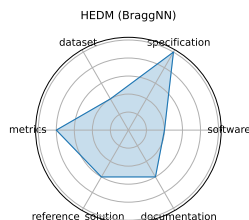
ratings.reference_solution.reason: BraggNN model is described and evaluated, but no direct implementation or inference scripts available

ratings.documentation.rating: 3

ratings.documentation.reason: Paper is clear, but lacks a GitHub repo or full reproducibility pipeline

id: hedm_braggnn

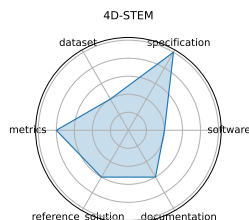
Citations: [39]



Ratings:

39 4D-STEM

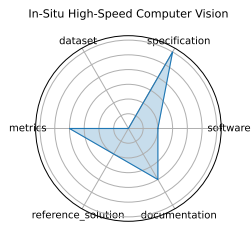
date: 2023-12-03
version: v1.0
last_updated: 2023-12
expired: unknown
valid: yes
valid_date: 2023-12-03
url: <https://openreview.net/pdf?id=7yt3N0o0W9>
doi: unknown
domain: Material Science
focus: Real-time ML for scanning transmission electron microscopy
keywords: - 4D-STEM - electron microscopy - real-time - image processing
summary: Proposes ML methods for real-time analysis of 4D scanning transmission electron microscopy datasets; framework details in progress.
licensing: unknown
task_types: - Image Classification - Streamed data inference
ai_capability_measured: - Real-time large-scale microscopy inference
metrics: - Classification accuracy - Throughput
models: - CNN models (prototype)
ml_motif: - Real-time, Image/CV
type: Model
ml_task: - Image Classification
solutions: 0
notes: In-progress; model design under development.
contact.name: Shuyu Qin
contact.email: shq219@lehigh.edu
results.links.name: ChatGPT LLM
fair.reproducible: in progress
fair.benchmark_ready: False
ratings.software.rating: 2
ratings.software.reason: No standalone code repository or setup instructions provided
ratings.specification.rating: 5
ratings.specification.reason: None
ratings.dataset.rating: 2
ratings.dataset.reason: No dataset links or FAIR metadata; unclear public access
ratings.metrics.rating: 4
ratings.metrics.reason: Only localization accuracy and inference time mentioned; not formally benchmarked with scripts
ratings.reference_solution.rating: 3
ratings.reference_solution.reason: BraggNN model is described and evaluated, but no direct implementation or inference scripts available
ratings.documentation.rating: 3
ratings.documentation.reason: Paper is clear, but lacks a GitHub repo or full reproducibility pipeline
id: d-stem
Citations: [40]



Ratings:

40 In-Situ High-Speed Computer Vision

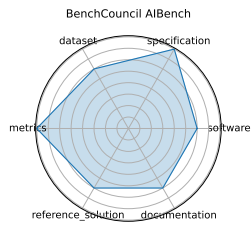
date: 2023-12-05
version: v1.0
last_updated: 2023-12
expired: unknown
valid: yes
valid_date: 2023-12-05
url: <https://arxiv.org/abs/2312.00128>
doi: 10.48550/arXiv.2312.00128
domain: Fusion/Plasma
focus: Real-time image classification for in-situ plasma diagnostics
keywords: - plasma - in-situ vision - real-time ML
summary: Applies low-latency CNN models for image classification of plasma diagnostics streams; supports deployment on embedded platforms.
licensing: Via Fermilab
task_types: - Image Classification
ai_capability_measured: - Real-time diagnostic inference
metrics: - Accuracy - FPS
models: - CNN
ml_motif: - Real-time, Image/CV
type: Model
ml_task: - Image Classification
solutions: Solution details are described in the referenced paper or repository.
notes: Embedded/deployment details in progress.
contact.name: unknown
contact.email: unknown
results.links.name: ChatGPT LLM
results.links.url: https://docs.google.com/document/d/1EqkRHuQs1yQqMvZs_L6p9JAY2vKX5OCTubzttFBuRoQ/edit?usp=sharing
fair.reproducible: in progress
fair.benchmark_ready: False
ratings.software.rating: 1
ratings.software.reason: No public implementation or containerized setup released
ratings.specification.rating: 3
ratings.specification.reason: No standardized I/O, latency constraint, or complete framing
ratings.dataset.rating: 0
ratings.dataset.reason: Dataset not provided or described in any formal way
ratings.metrics.rating: 2
ratings.metrics.reason: Throughput and accuracy mentioned, but not defined or benchmarked
ratings.reference_solution.rating: 1
ratings.reference_solution.reason: Prototype CNNs described; no code, baseline, or training details available
ratings.documentation.rating: 2
ratings.documentation.reason: Some insight via papers, but no working repo, setup, or replication path
id: in-situ_high-speed_computer_vision
Citations: [41]



Ratings:

41 BenchCouncil AIBench

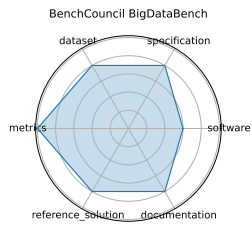
date: 2020-01-01
version: v1.0
last_updated: 2020-01
expired: unknown
valid: yes
valid_date: 2020-01-01
url: <https://www.benchcouncil.org/AIBench/>
doi: 10.48550/arXiv.1908.08998
domain: General
focus: End-to-end AI benchmarking across micro, component, and application levels
keywords: - benchmarking - AI systems - application-level evaluation
summary: AIBench is a comprehensive benchmark suite that evaluates AI workloads at different levels (micro, component, application) across hardware systems-covering image generation, object detection, translation, recommendation, video prediction, etc.
licensing: Apache License 2.0
task_types: - Training - Inference - End-to-end AI workloads
ai_capability_measured: - System-level AI workload performance
metrics: - Throughput - Latency - Accuracy
models: - ResNet - BERT - GANs - Recommendation systems
ml_motif: - General
type: Benchmark
ml_task: - NA
solutions: Solution details are described in the referenced paper or repository.
notes: Covers scenario-distilling, micro, component, and end-to-end benchmarks.
contact.name: Wanling Gao (BenchCouncil)
contact.email: unknown
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 3
ratings.software.reason: No containerized or automated implementation provided for full benchmark suite
ratings.specification.rating: 4
ratings.specification.reason: Task coverage is broad and well-scoped, but system constraints and expected outputs are not uniformly defined
ratings.dataset.rating: 3
ratings.dataset.reason: Multiple datasets are mentioned, but not consistently FAIR-documented, versioned, or linked
ratings.metrics.rating: 4
ratings.metrics.reason: Metrics are appropriate, but standardization and reproducibility across tasks vary
ratings.reference_solution.rating: 3
ratings.reference_solution.reason: Reference models (e.g., ResNet, BERT) described; no turnkey implementation or results repository for all levels
ratings.documentation.rating: 3
ratings.documentation.reason: Paper is comprehensive, but minimal user-facing documentation or structured reproduction guide
id: benchcouncil_aibench
Citations: [42]



Ratings:

42 BenchCouncil BigDataBench

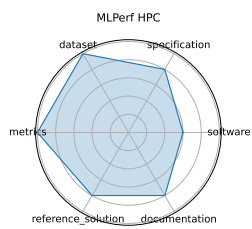
date: 2020-01-01
version: v1.0
last_updated: 2020-01
expired: unknown
valid: yes
valid_date: 2020-01-01
url: <https://www.benchcouncil.org/BigDataBench/>
doi: 10.48550/arXiv.1802.08254
domain: General
focus: Big data and AI benchmarking across structured, semi-structured, and unstructured data workloads
keywords: - big data - AI benchmarking - data analytics
summary: BigDataBench provides benchmarks for evaluating big data and AI workloads with realistic datasets (13 sources) and pipelines across analytics, graph, warehouse, NoSQL, streaming, and AI.
licensing: Apache License 2.0
task_types: - Data preprocessing - Inference - End-to-end data pipelines
ai_capability_measured: - Data processing and AI model inference performance at scale
metrics: - Data throughput - Latency - Accuracy
models: - CNN - LSTM - SVM - XGBoost
ml_motif: - General
type: Benchmark
ml_task: - NA
solutions: Solution details are described in the referenced paper or repository.
notes: Built on eight data motifs; provides Hadoop, Spark, Flink, MPI implementations.
contact.name: Jianfeng Zhan (BenchCouncil)
contact.email: unknown
results.links.name: ChatGPT LLM
results.links.url: <https://docs.google.com/document/d/1VFRxhR2G5A83S8PqKBrP99LLVgcCGvX2WW4vTtwxmQ4/edit?usp=sharing>
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 3
ratings.software.reason: No automated setup across all tasks; some components require manual integration.
ratings.specification.rating: 4
ratings.specification.reason: Specific I/O formats and hardware constraints are not uniformly detailed across all tasks.
ratings.dataset.rating: 4
ratings.dataset.reason: Some datasets lack consistent versioning or rich metadata annotations.
ratings.metrics.rating: 5
ratings.metrics.reason: None
ratings.reference_solution.rating: 4
ratings.reference_solution.reason: Not all benchmark components have fully reproducible baselines; deployment across platforms is fragmented.
ratings.documentation.rating: 4
ratings.documentation.reason: Setup requires manual steps; some task-specific instructions lack clarity.
id: benchcouncil_bigdatabench
Citations: [43]



Ratings:

43 MLPerf HPC

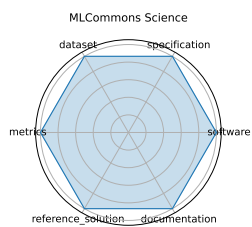
date: 2021-10-20
version: v1.0
last_updated: 2021-10
expired: unknown
valid: yes
valid_date: 2021-10-20
url: <https://github.com/mlcommons/hpc>
doi: 10.48550/arXiv.2110.11466
domain: Cosmology, Climate, Protein Structure, Catalysis
focus: Scientific ML training and inference on HPC systems
keywords: - HPC - training - inference - scientific ML
summary: MLPerf HPC introduces scientific model benchmarks (e.g., CosmoFlow, DeepCAM) aimed at large-scale HPC evaluation with >10x performance scaling through system-level optimizations.
licensing: Apache License 2.0
task_types: - Training - Inference
ai_capability_measured: - Scaling efficiency - training time - model accuracy on HPC
metrics: - Training time - Accuracy - GPU utilization
models: - CosmoFlow - DeepCAM - OpenCatalyst
ml_motif: - HPC/inference, HPC/training
type: Framework
ml_task: - NA
solutions: Solution details are described in the referenced paper or repository.
notes: Shared framework with MLCommons Science; reference implementations included.
contact.name: Steven Farrell (MLCommons)
contact.email: unknown
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 3
ratings.software.reason: Reference implementations exist but containerization and environment setup require manual effort across HPC systems.
ratings.specification.rating: 4
ratings.specification.reason: Hardware constraints and I/O formats are not fully defined for all scenarios.
ratings.dataset.rating: 5
ratings.dataset.reason: Not all data is independently versioned or comes with standardized FAIR metadata.
ratings.metrics.rating: 5
ratings.metrics.reason: None
ratings.reference_solution.rating: 4
ratings.reference_solution.reason: Reproducibility and environment tuning depend on system configuration; baseline models not uniformly bundled.
ratings.documentation.rating: 4
ratings.documentation.reason: Central guidance is available but requires domain-specific effort to replicate results across systems.
id: mlperf_hpc
Citations: [44]



Ratings:

44 MLCommons Science

date: 2023-06-01
version: v1.0
last_updated: 2023-06
expired: unknown
valid: yes
valid_date: 2023-06-01
url: <https://github.com/mlcommons/science>
doi: unknown
domain: Earthquake, Satellite Image, Drug Discovery, Electron Microscope, CFD
focus: AI benchmarks for scientific applications including time-series, imaging, and simulation
keywords: - science AI - benchmark - MLCommons - HPC
summary: MLCommons Science assembles benchmark tasks with datasets, targets, and implementations across earthquake forecasting, satellite imagery, drug screening, electron microscopy, and CFD to drive scientific ML reproducibility.
licensing: Apache License 2.0
task_types: - Time-series analysis - Image classification - Simulation surrogate modeling
ai_capability_measured: - Inference accuracy - simulation speed-up - generalization
metrics: - MAE - Accuracy - Speedup vs simulation
models: - CNN - GNN - Transformer
ml_motif: - Time-series, Image/CV, HPC/inference
type: Framework
ml_task: - NA
solutions: 0
notes: Joint national-lab effort under Apache-2.0 license.
contact.name: MLCommons Science Working Group
contact.email: unknown
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 5
ratings.software.reason: Actively maintained GitHub repository available at <https://github.com/mlcommons/science> with implementations, scripts, and reproducibility support.
ratings.specification.rating: 5
ratings.specification.reason: All five specification aspects are covered: system constraints, task, dataset format, benchmark inputs, and outputs.
ratings.dataset.rating: 5
ratings.dataset.reason: Public scientific datasets are used with defined splits. At least 4 FAIR principles are followed.
ratings.metrics.rating: 5
ratings.metrics.reason: Clearly defined metrics such as accuracy, training time, and GPU utilization are used. These metrics are explained and effectively capture solution performance.
ratings.reference_solution.rating: 5
ratings.reference_solution.reason: A reference implementation is available, well-documented, trainable/open, and includes full metric evaluation and software/hardware details.
ratings.documentation.rating: 5
ratings.documentation.reason: Thorough documentation exists covering the task, background, motivation, evaluation criteria, and includes a supporting paper.
id: mlcommons_science
Citations: [45]



Ratings:

45 LHC New Physics Dataset

date: 2021-07-05

version: v1.0

last_updated: 2021-07

expired: unknown

valid: yes

valid_date: 2021-07-05

url: <https://arxiv.org/pdf/2107.02157>

doi: unknown

domain: Particle Physics; Real-time Triggering

focus: Real-time LHC event filtering for anomaly detection using proton collision data

keywords: - anomaly detection - proton collision - real-time inference - event filtering - unsupervised ML

summary: A dataset of proton-proton collision events emulating a 40 MHz real-time data stream from LHC detectors, pre-filtered on electron or muon presence. Designed for unsupervised new-physics detection algorithms under latency/bandwidth constraints.

licensing: unknown

task_types: - Anomaly detection - Event classification

ai_capability_measured: - Unsupervised signal detection under latency and bandwidth constraints

metrics: - ROC-AUC - Detection efficiency

models: - Autoencoder - Variational autoencoder - Isolation forest

ml_motif: - Multiple

type: Framework

ml_task: - NA

solutions: 0

notes: Includes electron/muon-filtered background and black-box signal benchmarks; 1M events per black box.

contact.name: Ema Puljak (ema.puljak@cern.ch)

contact.email: unknown

datasets.links.name: Zenodo stores, background + 3 black-box signal sets. 1M events each

results.links.name: ChatGPT LLM

fair.reproducible: Yes

fair.benchmark_ready: Yes

ratings.software.rating: 3

ratings.software.reason: While not formally evaluated in the previous version, Zenodo and paper links suggest available code for baseline models (e.g., autoencoders, GANs), though they are scattered and not unified in a single repository.

ratings.specification.rating: 3

ratings.specification.reason: The task and context are clearly described, but system constraints and formal inputs/outputs are not fully specified.

ratings.dataset.rating: 5

ratings.dataset.reason: Large-scale dataset hosted on Zenodo, publicly available, well-documented, with defined train/test structure. Appears to follow at least 4 FAIR principles.

ratings.metrics.rating: 4

ratings.metrics.reason: Uses reasonable metrics (ROC-AUC, detection efficiency) that capture performance but lacks full explanation and standard evaluation tools.

ratings.reference_solution.rating: 2

ratings.reference_solution.reason: Baselines are described across multiple papers but lack centralized, reproducible implementations and hardware/software setup details.

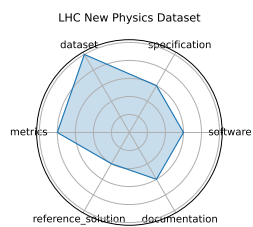
ratings.documentation.rating: 3

ratings.documentation.reason: Some description in papers and dataset metadata exists, but lacks a unified guide, README, or training setup in a central location.

id: lhc_new_physics_dataset

Citations: [46]

Ratings:



46 MLCommons Medical AI

date: 2023-07-17

version: v1.0

last_updated: 2023-07

expired: unknown

valid: yes

valid_date: 2023-07-17

url: <https://github.com/mlcommons/medical>

doi: unknown

domain: Healthcare; Medical AI

focus: Federated benchmarking and evaluation of medical AI models across diverse real-world clinical data

keywords: - medical AI - federated evaluation - privacy-preserving - fairness - healthcare benchmarks

summary: The MLCommons Medical AI working group develops benchmarks, best practices, and platforms (MedPerf, GaNDLF, COFE) to accelerate robust, privacy-preserving AI development for healthcare. MedPerf enables federated testing of clinical models on diverse datasets, improving generalizability and equity while keeping data onsite .

licensing: Apache License 2.0

task_types: - Federated evaluation - Model validation

ai_capability_measured: - Clinical accuracy - fairness - generalizability - privacy compliance

metrics: - ROC AUC - Accuracy - Fairness metrics

models: - MedPerf-validated CNNs - GaNDLF workflows

ml_motif: - Multiple

type: Platform

ml_task: - NA

solutions: 0

notes: Open-source platform under Apache-2.0; used across 20+ institutions and hospitals .

contact.name: Alex Karargyris (MLCommons Medical AI)

contact.email: unknown

datasets.links.name: Multi-institutional clinical datasets, radiology

results.links.name: ChatGPT LLM

fair.reproducible: Yes

fair.benchmark_ready: Yes

ratings.software.rating: 5

ratings.software.reason: GitHub repository (<https://github.com/mlcommons/medical>) provides actively maintained open-source tools like MedPerf and GaNDLF for federated medical AI evaluation.

ratings.specification.rating: 4

ratings.specification.reason: The platform defines federated tasks and model evaluation scenarios. Some clinical and system-level constraints are implied but not uniformly formalized across all use cases.

ratings.dataset.rating: 4

ratings.dataset.reason: Multi-institutional datasets used in federated settings; real-world data is handled privately onsite, but some FAIR aspects (e.g., accessibility and metadata) are implicit.

ratings.metrics.rating: 5

ratings.metrics.reason: Metrics such as ROC AUC, accuracy, and fairness are clearly specified and directly support goals like generalizability and equity.

ratings.reference_solution.rating: 3

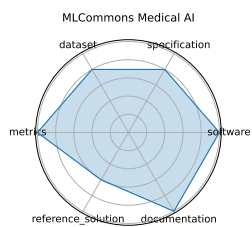
ratings.reference_solution.reason: GaNDLF workflows and MedPerf-validated CNNs are referenced, but not all baseline models are centrally documented or easily reproducible.

ratings.documentation.rating: 5

ratings.documentation.reason: Extensive documentation, papers, and community support exist. Clear examples and usage instructions are provided in GitHub and publications.

id: mlcommons_medical_ai

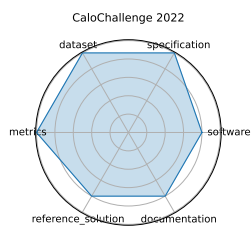
Citations: [47]



Ratings:

47 CaloChallenge 2022

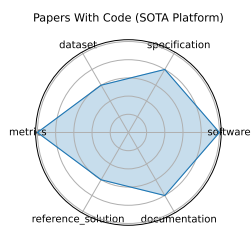
date: 2024-10-28
version: v1.0
last_updated: 2024-10
expired: unknown
valid: yes
valid_date: 2024-10-28
url: <http://arxiv.org/abs/2410.21611>
doi: 10.48550/arXiv.2410.21611
domain: LHC Calorimeter; Particle Physics
focus: Fast generative-model-based calorimeter shower simulation evaluation
keywords: - calorimeter simulation - generative models - surrogate modeling - LHC - fast simulation
summary: The Fast Calorimeter Simulation Challenge 2022 assessed 31 generative-model submissions (VAEs, GANs, Flows, Diffusion) on four calorimeter shower datasets; benchmarking shower quality, generation speed, and model complexity .
licensing: Via Fermilab
task_types: - Surrogate modeling
ai_capability_measured: - Simulation fidelity - speed - efficiency
metrics: - Histogram similarity - Classifier AUC - Generation latency
models: - VAE variants - GAN variants - Normalizing flows - Diffusion models
ml_motif: - Surrogate
type: Dataset
ml_task: - Surrogate Modeling
solutions: Solution details are described in the referenced paper or repository.
notes: The most comprehensive survey to date on ML-based calorimeter simulation; 31 submissions over different dataset sizes.
contact.name: Claudius Krause (CaloChallenge Lead)
contact.email: unknown
datasets.links.name: Four LHC calorimeter shower datasets
datasets.links.url: various voxel resolutions
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 4
ratings.software.reason: Community GitHub repos and model implementations are available for the 31 submissions. While not fully unified in one place, the software is accessible and reproducible.
ratings.specification.rating: 5
ratings.specification.reason: The task—evaluating fast generative calorimeter simulations—is clearly defined with benchmarking protocols, constraints like latency and model complexity, and structured evaluation criteria.
ratings.dataset.rating: 5
ratings.dataset.reason: Four well-structured calorimeter datasets are provided, with different voxel resolutions, open access, signal/background separation, and metadata. FAIR principles are well covered.
ratings.metrics.rating: 5
ratings.metrics.reason: Metrics like histogram similarity, classifier AUC, and generation latency are well defined and relevant for simulation quality, fidelity, and performance.
ratings.reference_solution.rating: 4
ratings.reference_solution.reason: Several baselines (GANs, VAEs, flows, diffusion models) are documented and evaluated. Some are available via community repos, though not all are fully standardized or bundled.
ratings.documentation.rating: 4
ratings.documentation.reason: Accompanied by a detailed paper and dataset description. Reproduction of pipelines may require additional setup or familiarity with the model submissions.
id: calochallenge_
Citations: [48]



Ratings:

48 Papers With Code (SOTA Platform)

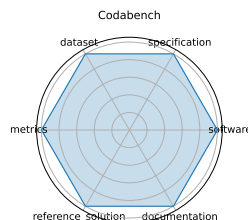
date: ongoing
version: v1.0
last_updated: 2025-06
expired: unknown
valid: yes
valid_date: ongoing
url: <https://paperswithcode.com/sota>
doi: unknown
domain: General ML; All domains
focus: Open platform tracking state-of-the-art results, benchmarks, and implementations across ML tasks and papers
keywords: - leaderboard - benchmarking - reproducibility - open-source
summary: Papers With Code (PWC) aggregates benchmark suites, tasks, and code across ML research: 12,423 benchmarks, 5,358 unique tasks, and 154,766 papers with code links. It tracks SOTA metrics and fosters reproducibility.
licensing: Apache License 2.0
task_types: - Multiple (Classification, Detection, NLP, etc.)
ai_capability_measured: - Model performance across tasks (accuracy - F1 - BLEU - etc.)
metrics: - Task-specific (Accuracy, F1, BLEU, etc.)
models: - All published models with code
ml_motif: - Multiple
type: Platform
ml_task: - Multiple
solutions: 0
notes: Community-driven open platform; automatic data extraction and versioning.
contact.name: Papers With Code Team
contact.email: unknown
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 5
ratings.software.reason: Actively maintained open-source platform (<https://paperswithcode.com>) under Apache 2.0 license; includes automatic integration with GitHub, datasets, and models for reproducibility.
ratings.specification.rating: 4
ratings.specification.reason: Task and benchmark structures are well organized and standardized, but due to its broad coverage, input/output formats vary significantly between tasks and are not always tightly controlled.
ratings.dataset.rating: 3
ratings.dataset.reason: Relies on external datasets submitted by the community. While links are available, FAIR compliance is not guaranteed or systematically enforced across all benchmarks.
ratings.metrics.rating: 5
ratings.metrics.reason: Tracks state-of-the-art using task-specific metrics like Accuracy, F1, BLEU, etc., with consistent aggregation and historical SOTA tracking.
ratings.reference_solution.rating: 3
ratings.reference_solution.reason: Provides links to implementations of many SOTA models, but no single unified reference baseline is required or maintained per benchmark.
ratings.documentation.rating: 4
ratings.documentation.reason: Strong front-end documentation and metadata on benchmarks, tasks, and models; however, some benchmark-specific instructions are sparse or dependent on external paper links.
id: papers_with_code_sota_platform
Citations: [49]



Ratings:

49 Codabench

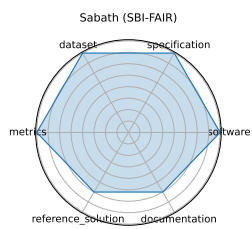
date: 2022-01-01
version: v1.0
last_updated: 2025-03
expired: unknown
valid: yes
valid_date: 2022-01-01
url: <https://www.codabench.org/>
doi: <https://doi.org/10.1016/j.patter.2022.100543>
domain: General ML; Multiple
focus: Open-source platform for organizing reproducible AI benchmarks and competitions
keywords: - benchmark platform - code submission - competitions - meta-benchmark
summary: Codabench (successor to CodaLab) is a flexible, easy-to-use, reproducible API platform for hosting AI benchmarks and code-submission challenges. It supports custom scoring, inverted benchmarks, and scalable public or private queues .
licensing: <https://github.com/codalab/codalab-competitions/wiki/Privacy>
task_types: - Multiple
ai_capability_measured: - Model reproducibility - performance across datasets
metrics: - Submission count - Leaderboard ranking - Task-specific metrics
models: - Arbitrary code submissions
ml_motif: - Multiple
type: Platform
ml_task: - Multiple
solutions: Several
notes: Hosts 51 public competitions, ~26 k users, 177 k submissions
contact.name: Isabelle Guyon (Université Paris-Saclay)
contact.email: unknown
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 1
ratings.software.reason: This is a platform for posting benchmarks, not a benchmark in itself.
ratings.specification.rating: 1
ratings.specification.reason: This is a platform for posting benchmarks, not a benchmark in itself.
ratings.dataset.rating: 1
ratings.dataset.reason: This is a platform for posting benchmarks, not a benchmark in itself.
ratings.metrics.rating: 1
ratings.metrics.reason: This is a platform for posting benchmarks, not a benchmark in itself.
ratings.reference_solution.rating: 1
ratings.reference_solution.reason: This is a platform for posting benchmarks, not a benchmark in itself.
ratings.documentation.rating: 1
ratings.documentation.reason: This is a platform for posting benchmarks, not a benchmark in itself.
id: codabench
Citations: [50]



Ratings:

50 Sabath (SBI-FAIR)

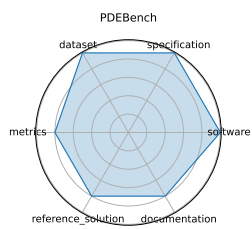
date: 2021-09-27
version: v1.0
last_updated: 2023-07
expired: unknown
valid: yes
valid_date: 2021-09-27
url: <https://sbi-fair.github.io/docs/software/sabath/>
doi: unknown
domain: Systems; Metadata
focus: FAIR metadata framework for ML-driven surrogate workflows in HPC systems
keywords: - meta-benchmark - metadata - HPC - surrogate modeling
summary: Sabath is a metadata framework from the SBI-FAIR group (UTK, Argonne, Virginia) facilitating FAIR-compliant benchmarking and surrogate execution logging across HPC systems .
licensing: BSD 3-Clause License
task_types: - Systems benchmarking
ai_capability_measured: - Metadata tracking - reproducible HPC workflows
metrics: - Metadata completeness - FAIR compliance
models: - NA
ml_motif: - Systems
type: Platform
ml_task: - NA
solutions: 0
notes: Developed by PI Piotr Luszczyk at UTK; integrates with MiniWeatherML, AutoPhaseNN, Cosmoflow, etc.
contact.name: Piotr Luszczyk
contact.email: luszczyk@utk.edu
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: N/A
ratings.software.rating: 4
ratings.software.reason: Actively maintained GitHub repository (<https://github.com/icl-utk-edu/slip/tree/sabath>) with BSD-licensed tooling for FAIR metadata capture; integrates with existing surrogate modeling benchmarks.
ratings.specification.rating: 4
ratings.specification.reason: FAIR metadata structure and logging goals are clearly described. Input/output definitions are implied through integrations (e.g., MiniWeatherML), though not always formalized.
ratings.dataset.rating: 4
ratings.dataset.reason: Datasets used in surrogate benchmarks are publicly available, well-structured, and FAIR-aligned, but not independently hosted by Sabath itself.
ratings.metrics.rating: 4
ratings.metrics.reason: Emphasizes metadata completeness and FAIR compliance. Metrics are clear and well-matched to its metadata-focused benchmarking context.
ratings.reference_solution.rating: 3
ratings.reference_solution.reason: Includes integration with multiple surrogate benchmarks and models, though not all are fully documented or packaged as standardized reference solutions.
ratings.documentation.rating: 3
ratings.documentation.reason: Basic instructions and code are provided on GitHub, but more detailed walkthroughs, use-case examples, or tutorials are limited.
id: sabath_sbi-fair
Citations: [51]



Ratings:

51 PDEBench

date: 2022-10-13
version: v0.1.0
last_updated: 2025-05
expired: unknown
valid: yes
valid_date: 2022-10-13
url: <https://github.com/pdebench/PDEBench>
doi: 10.48550/arXiv.2210.07182
domain: CFD; Weather Modeling
focus: Benchmark suite for ML-based surrogates solving time-dependent PDEs
keywords: - PDEs - CFD - scientific ML - surrogate modeling - NeurIPS
summary: PDEBench offers forward/inverse PDE tasks with large ready-to-use datasets and baselines (FNO, U-Net, PINN), packaged via a unified API. It won the SimTech Best Paper Award 2023 .
licensing: Other
task_types: - Supervised Learning
ai_capability_measured: - Time-dependent PDE modeling; physical accuracy
metrics: - RMSE - boundary RMSE - Fourier RMSE
models: - FNO - U-Net - PINN - Gradient-Based inverse methods
ml_motif: - Multiple
type: Framework
ml_task: - Supervised Learning
solutions: Solution details are described in the referenced paper or repository.
notes: Datasets hosted on DaRUS (DOI:10.18419/darus-2986); contact maintainers by email
contact.name: Makoto Takamoto (makoto.takamoto@neclab.eu)
contact.email: unknown
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 5
ratings.software.reason: GitHub repository (<https://github.com/pdebench/PDEBench>) is actively maintained and includes training pipelines, data loaders, and evaluation scripts. Installation and usage are well-documented.
ratings.specification.rating: 5
ratings.specification.reason: Clearly defined tasks for forward and inverse PDE problems, with structured input/output formats, system constraints, and task specifications.
ratings.dataset.rating: 5
ratings.dataset.reason: Diverse PDE datasets (synthetic and real-world) hosted on DaRUS with DOIs. Datasets are well-documented, structured, and follow FAIR practices.
ratings.metrics.rating: 4
ratings.metrics.reason: Includes RMSE, boundary RMSE, and Fourier-domain RMSE. These are well-suited to PDE problems, though rationale behind metric choices could be expanded in some cases.
ratings.reference_solution.rating: 4
ratings.reference_solution.reason: Baselines (FNO, U-Net, PINN, etc.) are available and documented, but not every model includes full training and evaluation reproducibility out-of-the-box.
ratings.documentation.rating: 4
ratings.documentation.reason: Strong documentation on GitHub including examples, configs, and usage instructions. Some model-specific details and tutorials could be further expanded.
id: pdebench
Citations: [52]



Ratings:

52 The Well

date: 2024-12-03

version: v1.0

last_updated: 2025-06

expired: unknown

valid: yes

valid_date: 2024-12-03

url: https://polymathic-ai.org/the_well/

doi: unknown

domain: biological systems, fluid dynamics, acoustic scattering, astrophysical MHD

focus: Foundation model + surrogate dataset spanning 16 physical simulation domains

keywords: - surrogate modeling - foundation model - physics simulations - spatiotemporal dynamics

summary: A 15 TB collection of ML-ready physics simulation datasets (HDF5), covering 16 domains-from biology to astro-physical magnetohydrodynamic simulations-with unified API and metadata. Ideal for training surrogate and foundation models on scientific data.

licensing: BSD 3-Clause License

task_types: - Supervised Learning

ai_capability_measured: - Surrogate modeling - physics-based prediction

metrics: - Dataset size - Domain breadth

models: - FNO baselines - U-Net baselines

ml_motif: - Foundation model, Surrogate

type: Dataset

ml_task: - Supervised Learning

solutions: 1

notes: Includes unified API and dataset metadata; see 2025 NeurIPS paper for full benchmark details. Size: 15 TB.

contact.name: Ruben Ohana

contact.email: rohana@flatironinstitute.org

datasets.links.name: 16 simulation datasets

datasets.links.url: HDF5) via PyPI/GitHub

results.links.name: ChatGPT LLM

fair.reproducible: Yes

fair.benchmark_ready: Yes

ratings.software.rating: 5

ratings.software.reason: BSD-licensed software and unified API are available via GitHub and PyPI. Supports loading and manipulating large HDF5 datasets across 16 domains.

ratings.specification.rating: 4

ratings.specification.reason: The benchmark includes clearly defined surrogate modeling tasks, data structure, and meta-data. However, constraints and formal task specs vary slightly across domains.

ratings.dataset.rating: 5

ratings.dataset.reason: 15 TB of ML-ready HDF5 datasets across 16 physics domains. Public, well-structured, richly annotated, and designed with FAIR principles in mind.

ratings.metrics.rating: 3

ratings.metrics.reason: Domain breadth and dataset size are emphasized. Standardized quantitative metrics for model evaluation (e.g., RMSE, accuracy) are not uniformly applied across all domains.

ratings.reference_solution.rating: 3

ratings.reference_solution.reason: Includes FNO and U-Net baselines, but does not yet provide fully trained, reproducible models or scripts across all datasets.

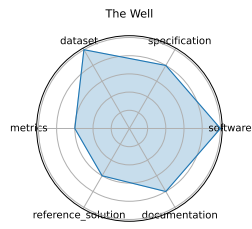
ratings.documentation.rating: 4

ratings.documentation.reason: The GitHub repo and NeurIPS paper provide detailed guidance on dataset use, structure, and training setup. Tutorials and walkthroughs could be expanded further.

id: the_well

Citations: [53]

Ratings:



53 LLM-Inference-Bench

date: 2024-10-31

version: v1.0

last_updated: 2024-11

expired: unknown

valid: yes

valid_date: 2024-10-31

url: <https://github.com/argonne-lcf/LLM-Inference-Bench>

doi: unknown

domain: LLM; HPC/inference

focus: Hardware performance benchmarking of LLMs on AI accelerators

keywords: - LLM - inference benchmarking - GPU - accelerator - throughput

summary: A suite evaluating inference performance of LLMs (LLaMA, Mistral, Qwen) across diverse accelerators (NVIDIA, AMD, Intel, SambaNova) and frameworks (vLLM, DeepSpeed-MII, etc.), with an interactive dashboard and per-platform metrics.

licensing: BSD 3-Clause "New" or "Revised" License

task_types: - Inference Benchmarking

ai_capability_measured: - Inference throughput - latency - hardware utilization

metrics: - Token throughput (tok/s) - Latency - Framework-hardware mix performance

models: - LLaMA-2-7B - LLaMA-2-70B - Mistral-7B - Qwen-7B

ml_motif: - HPC/inference

type: Dataset

ml_task: - Inference Benchmarking

solutions: 0

notes: Licensed under BSD-3, maintained by Argonne; supports GPUs and accelerators.

contact.name: Krishna Teja Chitty-Venkata (Argonne LCF)

contact.email: unknown

results.links.name: ChatGPT LLM

fair.reproducible: Yes

fair.benchmark_ready: Yes

ratings.software.rating: 5

ratings.software.reason: Public GitHub repository (<https://github.com/argonne-lcf/LLM-Inference-Bench>) under BSD-3 license. Includes scripts, configurations, and dashboards for running and visualizing LLM inference benchmarks across multiple accelerator platforms.

ratings.specification.rating: 5

ratings.specification.reason: Benchmark scope, models, accelerator targets, and supported frameworks are clearly specified. Input configurations and output metrics are standardized across hardware types.

ratings.dataset.rating: 2

ratings.dataset.reason: No novel dataset is introduced; benchmark relies on pre-trained LLMs and synthetic inference inputs. Dataset structure and FAIR considerations are minimal.

ratings.metrics.rating: 5

ratings.metrics.reason: Hardware-specific metrics (token throughput, latency, utilization) are well-defined, consistently measured, and aggregated in dashboards.

ratings.reference_solution.rating: 3

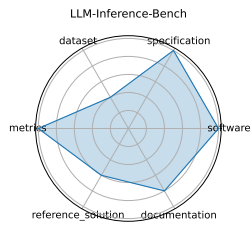
ratings.reference_solution.reason: Inference configurations and baseline performance results are provided, but there are no full reference training pipelines or model implementations.

ratings.documentation.rating: 4

ratings.documentation.reason: GitHub repo provides clear usage instructions, setup guides, and interactive dashboard tooling. Some areas like benchmarking extensions or advanced tuning are less detailed.

id: llm-inference-bench

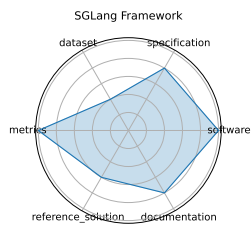
Citations: [54]



Ratings:

54 SGLang Framework

date: 2023-12-12
version: v0.4.9
last_updated: 2025-06
expired: unknown
valid: yes
valid_date: 2023-12-12
url: <https://github.com/sgl-project/sglang/tree/main/benchmark>
doi: 10.48550/arXiv.2312.07104
domain: LLM Vision
focus: Fast serving framework for LLMs and vision-language models
keywords: - LLM serving - vision-language - RadixAttention - performance - JSON decoding
summary: A high-performance open-source serving framework combining efficient backend runtime (RadixAttention, batching, quantization) and expressive frontend language, boosting LLM/VLM inference throughput up to ~3x over alternatives.
licensing: Apache License 2.0
task_types: - Model serving framework
ai_capability_measured: - Serving throughput - JSON/task-specific latency
metrics: - Tokens/sec - Time-to-first-token - Throughput gain vs baseline
models: - LLaVA - DeepSeek - Llama
ml_motif: - LLM Vision
type: Framework
ml_task: - Model serving
solutions: Solution details are described in the referenced paper or repository.
notes: Deployed in production (xAI, NVIDIA, Google Cloud); v0.4.8 release June 2025.
contact.name: SGLang Team
contact.email: unknown
datasets.links.name: Benchmark configs
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 5
ratings.software.reason: Actively maintained and production-deployed (e.g., xAI, NVIDIA); source code available under Apache 2.0. Includes efficient backends (RadixAttention, quantization, batching) and full serving infrastructure.
ratings.specification.rating: 4
ratings.specification.reason: The framework clearly defines performance targets, serving logic, and model integration. Input/output expectations are consistent, but not all benchmarks are standardized.
ratings.dataset.rating: 2
ratings.dataset.reason: Does not introduce new datasets; instead, it evaluates performance using existing model benchmarks. Only configuration files are included.
ratings.metrics.rating: 5
ratings.metrics.reason: Serving-related metrics such as tokens/sec, time-to-first-token, and throughput gain vs. baselines are well-defined and consistently applied.
ratings.reference_solution.rating: 3
ratings.reference_solution.reason: Provides benchmark configs and example integrations (e.g., with LLaVA, DeepSeek), but not all models or scripts are runnable out-of-the-box.
ratings.documentation.rating: 4
ratings.documentation.reason: Strong GitHub documentation, install guides, and benchmarks. Some advanced topics (e.g., scaling, hardware tuning) could use deeper walkthroughs.
id: sglang_framework
Citations: [55]



Ratings:

55 vLLM Inference and Serving Engine

date: 2023-09-12

version: v0.10.0

last_updated: 2025-06

expired: unknown

valid: yes

valid_date: 2023-09-12

url: <https://github.com/vllm-project/vllm/tree/main/benchmarks>

doi: unknown

domain: LLM; HPC/inference

focus: High-throughput, memory-efficient inference and serving engine for LLMs

keywords: - LLM inference - PagedAttention - CUDA graph - streaming API - quantization

summary: vLLM is a fast, high-throughput, memory-efficient inference and serving engine for large language models, featuring PagedAttention, continuous batching, and support for quantized and pipelined model execution. Benchmarks compare it to TensorRT-LLM, SGLang, and others.

licensing: Apache License 2.0

task_types: - Inference Benchmarking

ai_capability_measured: - Throughput - latency - memory efficiency

metrics: - Tokens/sec - Time to First Token (TTFT) - Memory footprint

models: - LLaMA - Mixtral - FlashAttention-based models

ml_motif: - HPC/inference

type: Framework

ml_task: - Inference

solutions: 0

notes: Incubated by LF AI and Data; achieves up to 24x throughput over HuggingFace Transformers

contact.name: Woosuk Kwon (vLLM Team)

contact.email: unknown

results.links.name: ChatGPT LLM

fair.reproducible: Yes

fair.benchmark_ready: Yes

ratings.software.rating: 5

ratings.software.reason: Actively maintained open-source project under Apache 2.0. GitHub repo includes full serving engine, benchmarking scripts, CUDA integration, and deployment examples.

ratings.specification.rating: 5

ratings.specification.reason: Inference benchmarks are well-defined with clear input/output formats and platform-specific constraints. Covers multiple models, hardware backends, and batching configurations.

ratings.dataset.rating: 3

ratings.dataset.reason: No traditional dataset is included. Instead, it uses structured configs and logs suitable for inference benchmarking. FAIR principles are only partially applicable.

ratings.metrics.rating: 5

ratings.metrics.reason: Comprehensive performance metrics like tokens/sec, time-to-first-token (TTFT), and memory footprint are consistently applied and benchmarked across frameworks.

ratings.reference_solution.rating: 4

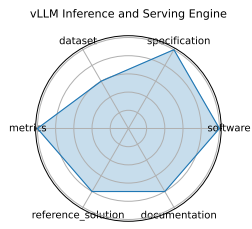
ratings.reference_solution.reason: Provides runnable scripts and configs for several models (LLaMA, Mixtral, etc.) across platforms. Baselines are reproducible, though not all models are fully wrapped or hosted.

ratings.documentation.rating: 4

ratings.documentation.reason: Well-structured GitHub documentation with setup instructions, config examples, benchmarking comparisons, and performance tuning guides.

id: vllm_inference_and_serving_engine

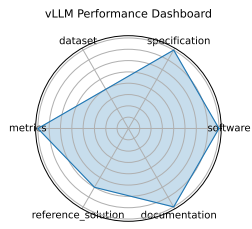
Citations: [56]



Ratings:

56 vLLM Performance Dashboard

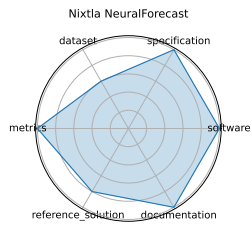
date: 2022-06-22
version: v1.0
last_updated: 2025-01
expired: unknown
valid: yes
valid_date: 2022-06-22
url: <https://simon-mo-workspace.observablehq.cloud/vllm-dashboard-v0/>
doi: unknown
domain: LLM; HPC/inference
focus: Interactive dashboard showing inference performance of vLLM
keywords: - Dashboard - Throughput visualization - Latency analysis - Metric tracking
summary: A live visual dashboard for vLLM showcasing throughput, latency, and other inference metrics across models and hardware configurations.
licensing: unknown
task_types: - Performance visualization
ai_capability_measured: - Throughput - latency - hardware utilization
metrics: - Tokens/sec - TTFT - Memory usage
models: - LLaMA-2 - Mistral - Qwen
ml_motif: - HPC/inference
type: Framework
ml_task: - Visualization
solutions: 0
notes: Built using ObservableHQ; integrates live data from vLLM benchmarks. The URL requires a login to access the content.
contact.name: Simon Mo
contact.email: unknown
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 4
ratings.software.reason: Interactive dashboard built with ObservableHQ and linked to vLLM benchmarks. Source code is not fully open, but backend integration with vLLM is well-maintained.
ratings.specification.rating: 4
ratings.specification.reason: While primarily a visualization tool, it includes benchmark configurations, metric definitions, and supports comparison across models and hardware.
ratings.dataset.rating: 2
ratings.dataset.reason: No datasets are bundled; the dashboard visualizes metrics derived from model inference logs or external endpoints, not a formal dataset.
ratings.metrics.rating: 4
ratings.metrics.reason: Tracks tokens/sec, TTFT, memory usage, and platform comparisons. Metrics are clear but focused on visualization rather than statistical robustness.
ratings.reference_solution.rating: 3
ratings.reference_solution.reason: Dashboards include reproducible views of benchmarked models, but do not ship with runnable model code. Relies on external serving infrastructure.
ratings.documentation.rating: 4
ratings.documentation.reason: Public dashboard with instructions and tooltips; documentation is clear, though access is restricted (login required) and backend setup is opaque to users.
id: vllm_performance_dashboard
Citations: [57]



Ratings:

57 Nixtla NeuralForecast

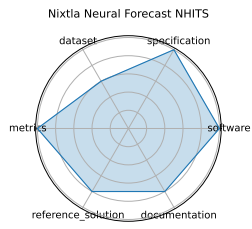
date: 2022-04-01
version: v3.0.2
last_updated: 2025-06
expired: unknown
valid: yes
valid_date: 2022-04-01
url: <https://github.com/Nixtla/neuralforecast>
doi: unknown
domain: Time-series forecasting; General ML
focus: High-performance neural forecasting library with >30 models
keywords: - time-series - neural forecasting - NBEATS, NHITS, TFT - probabilistic forecasting - usability
summary: NeuralForecast offers scalable, user-friendly implementations of over 30 neural forecasting models (NBEATS, NHITS, TFT, DeepAR, etc.), emphasizing quality, usability, interpretability, and performance.
licensing: Apache License 2.0
task_types: - Time-series forecasting
ai_capability_measured: - Forecast accuracy - interpretability - speed
metrics: - RMSE - MAPE - CRPS
models: - NBEATS - NHITS - TFT - DeepAR
ml_motif: - Time-series
type: Platform
ml_task: - Forecasting
solutions: 0
notes: AutoModel supports hyperparameter tuning and distributed execution via Ray and Optuna. First official NHITS implementation. contentReference oaicite:4 ndex=4
contact.name: Kin G. Olivares (Nixtla)
contact.email: unknown
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 5
ratings.software.reason: Actively maintained open-source library under Apache 2.0. Offers a clean API, extensive model zoo (>30 models), integration with Ray, Optuna, and supports scalable training and inference workflows.
ratings.specification.rating: 5
ratings.specification.reason: Forecasting task is well-defined with clear input/output structures. Framework supports probabilistic and deterministic forecasting, with unified interfaces and support for batch evaluation.
ratings.dataset.rating: 3
ratings.dataset.reason: NeuralForecast does not include its own datasets but supports standard datasets (e.g., M4, M5, ETT). FAIR compliance depends on user-supplied data.
ratings.metrics.rating: 5
ratings.metrics.reason: RMSE, MAPE, CRPS, and other domain-relevant metrics are well supported and integrated into the evaluation loop.
ratings.reference_solution.rating: 4
ratings.reference_solution.reason: Includes runnable model baselines and training scripts for all supported models. Some models have pretrained weights, but not all are fully benchmarked out-of-the-box.
ratings.documentation.rating: 5
ratings.documentation.reason: Rich documentation with examples, API references, tutorials, notebooks, and CLI support. PyPI, GitHub, and official blog posts offer clear guidance for usage and extension.
id: nixtla_neuralforecast
Citations: [58]



Ratings:

58 Nixtla Neural Forecast NHITS

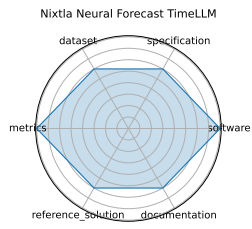
date: 2023-06-01
version: v3.0.2
last_updated: 2025-06
expired: unknown
valid: yes
valid_date: 2023-06-01
url: <https://github.com/Nixtla/neuralforecast>
doi: unknown
domain: Time-series; General ML
focus: Official NHITS implementation for long-horizon time series forecasting
keywords: - NHITS - long-horizon forecasting - neural interpolation - time-series
summary: NHITS (Neural Hierarchical Interpolation for Time Series) is a state-of-the-art model that improved accuracy by ~25% and reduced compute by 50x compared to Transformer baselines, using hierarchical interpolation and multi-rate sampling.
licensing: Apache License 2.0
task_types: - Time-series forecasting
ai_capability_measured: - Accuracy - compute efficiency for long series
metrics: - RMSE - MAPE
models: - NHITS
ml_motif: - Time-series
type: Platform
ml_task: - Forecasting
solutions: 0
notes: Official implementation in NeuralForecast, included since its AAAI 2023 release.
contact.name: Kin G. Olivares (Nixtla)
contact.email: unknown
datasets.links.name: Standard forecast datasets, M4
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 5
ratings.software.reason: Implemented within the open-source NeuralForecast library under Apache 2.0. Includes training, evaluation, and hyperparameter tuning pipelines. Actively maintained.
ratings.specification.rating: 5
ratings.specification.reason: The NHITS forecasting task is clearly defined with structured input/output formats. Model design targets long-horizon accuracy and compute efficiency.
ratings.dataset.rating: 3
ratings.dataset.reason: Uses standard benchmark datasets like M4, but does not bundle them directly. FAIR compliance depends on external dataset sources and user setup.
ratings.metrics.rating: 5
ratings.metrics.reason: Evaluated using RMSE, MAPE, and other standard forecasting metrics, integrated into training and evaluation APIs.
ratings.reference_solution.rating: 4
ratings.reference_solution.reason: Official NHITS implementation is fully reproducible with training/eval configs, though pretrained weights are not always provided.
ratings.documentation.rating: 4
ratings.documentation.reason: Well-documented on GitHub and in AAAI paper, with code examples, training guidance, and usage tutorials. More model-specific docs could improve clarity further.
id: nixtla_neural_forecast_nhits
Citations: [59]



Ratings:

59 Nixtla Neural Forecast TimeLLM

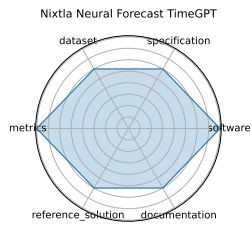
date: 2023-10-03
version: v3.0.2
last_updated: 2025-06
expired: unknown
valid: yes
valid_date: 2023-10-03
url: <https://github.com/Nixtla/neuralforecast>
doi: 10.48550/arXiv.2310.01728
domain: Time-series; General ML
focus: Reprogramming LLMs for time series forecasting
keywords: - Time-LLM - language model - time-series - reprogramming
summary: Time-LLM uses reprogramming layers to adapt frozen LLMs for time series forecasting, treating forecasting as a language task .
licensing: Apache License 2.0
task_types: - Time-series forecasting
ai_capability_measured: - Model reuse via LLM - few-shot forecasting
metrics: - RMSE - MAPE
models: - Time-LLM
ml_motif: - Time-series
type: Platform
ml_task: - Forecasting
solutions: Solution details are described in the referenced paper or repository.
notes: Fully open-source; transforms forecasting using LLM text reconstruction.
contact.name: Ming Jin (Nixtla)
contact.email: unknown
datasets.links.name: Standard forecast datasets, M4
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 4
ratings.software.reason: Fully open-source under Apache 2.0, integrated into the NeuralForecast library. Includes Time-LLM implementation with example usage and training scripts.
ratings.specification.rating: 3
ratings.specification.reason: High-level framing of forecasting as language modeling is clear, but detailed input/output specifications, constraints, and task formalization are minimal.
ratings.dataset.rating: 3
ratings.dataset.reason: Evaluated on standard datasets like M4 and ETT, but dataset splits and versioning are not bundled or explicitly FAIR-compliant.
ratings.metrics.rating: 4
ratings.metrics.reason: Standard forecasting metrics such as RMSE, MAPE, and SMAPE are reported. Evaluation is consistent, though deeper metric justification is limited.
ratings.reference_solution.rating: 3
ratings.reference_solution.reason: Time-LLM implementation is open and reproducible, but limited baselines or comparative implementations are included directly.
ratings.documentation.rating: 3
ratings.documentation.reason: GitHub README provides installation and quick usage examples, but lacks detailed API docs, training walkthroughs, or extended tutorials.
id: nixtla_neural_forecast_timellm
Citations: [60]



Ratings:

60 Nixtla Neural Forecast TimeGPT

date: 2023-10-05
version: v3.0.2
last_updated: 2025-06
expired: unknown
valid: yes
valid_date: 2023-10-05
url: <https://github.com/Nixtla/neuralforecast>
doi: 10.48550/arXiv.2310.03589
domain: Time-series; General ML
focus: Time-series foundation model "TimeGPT" for forecasting and anomaly detection
keywords: - TimeGPT - foundation model - time-series - generative model
summary: TimeGPT is a transformer-based generative pretrained model on 100B+ time series data for zero-shot forecasting and anomaly detection via API .
licensing: Apache License 2.0
task_types: - Time-series forecasting - Anomaly detection
ai_capability_measured: - Zero-shot forecasting - anomaly detection
metrics: - RMSE - Anomaly detection metrics
models: - TimeGPT
ml_motif: - Time-series
type: Platform
ml_task: - Forecasting
solutions: Solution details are described in the referenced paper or repository.
notes: Offered via Nixtla API and Azure Studio; enterprise-grade support available.
contact.name: Azul Garza (Nixtla)
contact.email: unknown
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 4
ratings.software.reason: Fully open-source Apache 2.0 implementation integrated in NeuralForecast, supporting training and evaluation via API. Production-grade deployment available via Nixtla API and Azure.
ratings.specification.rating: 3
ratings.specification.reason: Concept and forecasting goals are described, but formal input/output definitions and task constraints are not rigorously specified.
ratings.dataset.rating: 3
ratings.dataset.reason: Evaluated on existing open datasets, but consolidated data release, splits, and FAIR metadata are not provided.
ratings.metrics.rating: 4
ratings.metrics.reason: Uses standard forecasting metrics such as RMSE, MASE, SMAPE, and anomaly detection metrics consistently across evaluations.
ratings.reference_solution.rating: 3
ratings.reference_solution.reason: TimeGPT implementation is available, but baseline comparisons and additional reference models are limited.
ratings.documentation.rating: 3
ratings.documentation.reason: Basic README with installation and usage examples; more detailed API docs and tutorials would improve usability.
id: nixtla_neural_forecast_timegpt
Citations: [61]

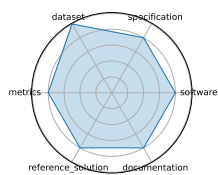


Ratings:

61 HDR ML Anomaly Challenge (Gravitational Waves)

date: 2025-03-03
version: v1.0
last_updated: 2025-03
expired: unknown
valid: yes
valid_date: 2025-03-03
url: <https://www.codabench.org/competitions/2626/>
doi: 10.48550/arXiv.2503.02112
domain: Astrophysics; Time-series
focus: Detecting anomalous gravitational-wave signals from LIGO/Virgo datasets
keywords: - anomaly detection - gravitational waves - astrophysics - time-series
summary: A benchmark for detecting anomalous transient gravitational-wave signals, including "unknown-unknowns," using preprocessed LIGO time-series at 4096 Hz. Competitors submit inference models on Codabench for continuous 50 ms segments from dual interferometers.
licensing: NA
task_types: - Anomaly detection
ai_capability_measured: - Novel event detection in physical signals
metrics: - ROC-AUC - Precision/Recall
models: - Deep latent CNNs - Autoencoders
ml_motif: - Time-series
type: Dataset
ml_task: - Anomaly detection
solutions: Solution details are described in the referenced paper or repository.
notes: NSF HDR A3D3 sponsored; prize pool and starter kit provided on Codabench.
contact.name: HDR A3D3 Team
contact.email: unknown
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 4
ratings.software.reason: Benchmark platform provided on Codabench with starter kits and submission infrastructure. Code and baseline models are publicly accessible but not extensively maintained beyond the challenge.
ratings.specification.rating: 4
ratings.specification.reason: Well-defined anomaly detection task on gravitational-wave time series with clear input/output expectations and challenge constraints.
ratings.dataset.rating: 5
ratings.dataset.reason: Uses preprocessed LIGO/Virgo time series data at 4096 Hz, publicly available and standard in astrophysics.
ratings.metrics.rating: 4
ratings.metrics.reason: ROC-AUC, precision, and recall metrics are clearly specified and appropriate for anomaly detection.
ratings.reference_solution.rating: 4
ratings.reference_solution.reason: Baseline deep latent CNNs and autoencoders are provided and reproducible, but not extensively documented.
ratings.documentation.rating: 4
ratings.documentation.reason: Documentation includes challenge instructions, starter kit details, and baseline descriptions, but could benefit from more thorough tutorials and code walkthroughs.
id: hdr_ml_anomaly_challenge_gravitational_waves
Citations: [62]

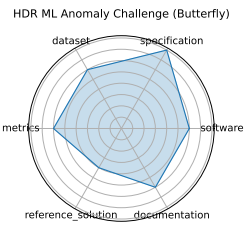
HDR ML Anomaly Challenge (Gravitational Waves)



Ratings:

62 HDR ML Anomaly Challenge (Butterfly)

date: 2025-03-03
version: v1.0
last_updated: 2025-03
expired: unknown
valid: yes
valid_date: 2025-03-03
url: <https://www.codabench.org/competitions/3764/>
doi: 10.48550/arXiv.2503.02112
domain: Genomics; Image/CV
focus: Detecting hybrid butterflies via image anomaly detection in genomic-informed dataset
keywords: - anomaly detection - computer vision - genomics - butterfly hybrids
summary: Image-based challenge for detecting butterfly hybrids in microscopy-driven species data. Participants evaluate models on Codabench using image segmentation/classification.
licensing: NA
task_types: - Anomaly detection
ai_capability_measured: - Hybrid detection in biological systems
metrics: - Classification accuracy - F1 score
models: - CNN-based detectors
ml_motif: - Image/CV
type: Dataset
ml_task: - Anomaly detection
solutions: Solution details are described in the referenced paper or repository.
notes: Hybrid detection benchmarks hosted on Codabench
contact.name: Imageomics/HDR Team
contact.email: unknown
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 3
ratings.software.reason: Codabench platform provides submission infrastructure but no fully maintained code repository or reproducible baseline implementations.
ratings.specification.rating: 4
ratings.specification.reason: Task is clearly described with domain-specific anomaly detection objectives and relevant physics motivation.
ratings.dataset.rating: 3
ratings.dataset.reason: Dataset consists of real detector data with synthetic anomaly injections; access is restricted and requires NDA, limiting openness and FAIR compliance.
ratings.metrics.rating: 3
ratings.metrics.reason: Standard metrics (ROC, F1, precision) are used; evaluation protocols are clear but not deeply elaborated.
ratings.reference_solution.rating: 2
ratings.reference_solution.reason: Baselines are partially described but lack public code or reproducible execution scripts.
ratings.documentation.rating: 3
ratings.documentation.reason: Challenge website provides basic descriptions and evaluation metrics but lacks comprehensive tutorials or example workflows.
id: hdr_ml_anomaly_challenge_butterfly
Citations: [63]

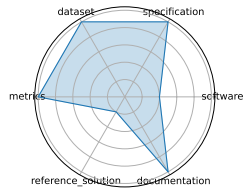


Ratings:

63 HDR ML Anomaly Challenge (Sea Level Rise)

date: 2025-03-03
version: v1.0
last_updated: 2025-03
expired: unknown
valid: yes
valid_date: 2025-03-03
url: <https://www.codabench.org/competitions/3223/>
doi: 10.48550/arXiv.2503.02112
domain: Climate Science; Time-series, Image/CV
focus: Detecting anomalous sea-level rise and flooding events via time-series and satellite imagery
keywords: - anomaly detection - climate science - sea-level rise - time-series - remote sensing
summary: A challenge combining North Atlantic sea-level time-series and satellite imagery to detect flooding anomalies. Models submitted via Codabench.
licensing: NA
task_types: - Anomaly detection
ai_capability_measured: - Detection of environmental anomalies
metrics: - ROC-AUC - Precision/Recall
models: - CNNs, RNNs, Transformers
ml_motif: - Time-series, Image/CV
type: Dataset
ml_task: - Anomaly detection
solutions: Solution details are described in the referenced paper or repository.
notes: Sponsored by NSF HDR; integrates sensor and satellite data.
contact.name: HDR A3D3 Team
contact.email: unknown
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 2
ratings.software.reason: Benchmark platform exists on Codabench, but no baseline code or maintained repository for reference solutions provided yet.
ratings.specification.rating: 5
ratings.specification.reason: Well-defined anomaly detection task combining satellite imagery and time-series data, with clear physical and domain-specific framing.
ratings.dataset.rating: 5
ratings.dataset.reason: Uses preprocessed, public, and well-structured sensor and satellite data for the North Atlantic sea-level rise region.
ratings.metrics.rating: 5
ratings.metrics.reason: Standard metrics such as ROC-AUC, precision, and recall are specified and suitable for the anomaly detection tasks.
ratings.reference_solution.rating: 1
ratings.reference_solution.reason: No starter models or baseline implementations linked or provided publicly.
ratings.documentation.rating: 5
ratings.documentation.reason: Challenge page, starter kits, and related papers offer strong guidance for participants.
id: hdr_ml_anomaly_challenge_sea_level_rise
Citations: [64]

HDR ML Anomaly Challenge (Sea Level Rise)

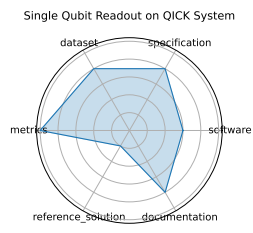


Ratings:

64 Single Qubit Readout on QICK System

date: 2025-01-24
version: v1.0
last_updated: 2025-02
expired: unknown
valid: yes
valid_date: 2025-01-24
url: <https://github.com/fastmachinelearning/ml-quantum-readout>
doi: 10.48550/arXiv.2501.14663
domain: Quantum Computing
focus: Real-time single-qubit state classification using FPGA firmware
keywords: - qubit readout - hls4ml - FPGA - QICK
summary: Implements real-time ML models for single-qubit readout on the Quantum Instrumentation Control Kit (QICK), using hls4ml to deploy quantized neural networks on RFSoc FPGAs. Offers high-fidelity, low-latency quantum state discrimination. :contentReference[oaicite:0]{index=0}
licensing: NA
task_types: - Classification
ai_capability_measured: - Single-shot fidelity - inference latency
metrics: - Accuracy - Latency
models: - hls4ml quantized NN
ml_motif: - Real-time
type: Benchmark
ml_task: - Supervised Learning
solutions: Solution details are described in the referenced paper or repository.
notes: Achieves ~96% fidelity with ~32 ns latency and low FPGA resource utilization.
contact.name: Javier Campos, Giuseppe Di Guglielmo
contact.email: unknown
datasets.links.name: Zenodo: ml-quantum-readout dataset
datasets.links.url: zenodo.org/records/14427490
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 3
ratings.software.reason: Code and FPGA firmware available on GitHub; integration with hls4ml demonstrated. Some deployment details and examples are provided but overall software maturity is moderate.
ratings.specification.rating: 4
ratings.specification.reason: Task clearly defined: real-time single-qubit state classification with latency and fidelity constraints. Labeling and ground truth definitions could be more explicit.
ratings.dataset.rating: 4
ratings.dataset.reason: Dataset hosted on Zenodo with structured data; however, detailed documentation on image acquisition and labeling pipeline is limited.
ratings.metrics.rating: 5
ratings.metrics.reason: Standard classification metrics (accuracy, latency) are used and directly relevant to the quantum readout task.
ratings.reference_solution.rating: 1
ratings.reference_solution.reason: No baseline or starter models with runnable code are linked publicly.
ratings.documentation.rating: 4
ratings.documentation.reason: Codabench task page and GitHub repo provide descriptions and usage instructions, but detailed API or deployment tutorials are limited.
id: single_qubit_readout_on_qick_system
Citations: [65]

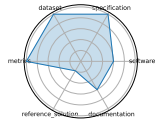
Ratings:



65 GPQA: A Graduate-Level Google-Proof Question and Answer Benchmark

date: 2023-11-20
version: v1.0
last_updated: 2023-11
expired: unknown
valid: yes
valid_date: 2023-11-20
url: <https://arxiv.org/abs/2311.12022>
doi: 10.48550/arXiv.2311.12022
domain: Science (Biology, Physics, Chemistry)
focus: Graduate-level, expert-validated multiple-choice questions hard even with web access
keywords: - Google-proof - multiple-choice - expert reasoning - science QA
summary: Contains 448 challenging questions written by domain experts, with expert accuracy at 65% (74% discounting clear errors) and non-experts reaching just 34%. GPT-4 baseline scores ~39%-designed for scalable oversight evaluation.
licensing: NA
task_types: - Multiple choice
ai_capability_measured: - Scientific reasoning - knowledge probing
metrics: - Accuracy
models: - GPT-4 baseline
ml_motif: - Multiple choice
type: Benchmark
ml_task: - Multiple choice
solutions: Solution details are described in the referenced paper or repository.
notes: Google-proof, supports oversight research.
contact.name: David Rein (NYU)
contact.email: unknown
datasets.links.name: GPQA dataset
datasets.links.url: [zip/HuggingFace](#)
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 3
ratings.software.reason: Dataset and benchmark materials are publicly available via HuggingFace and GitHub, but no integrated runnable code or software framework is provided.
ratings.specification.rating: 5
ratings.specification.reason: Task is clearly defined as a multiple-choice benchmark requiring expert-level scientific reasoning. Input/output formats and evaluation criteria are well described.
ratings.dataset.rating: 5
ratings.dataset.reason: The GPQA dataset is publicly released, well curated, with metadata and clearly documented splits.
ratings.metrics.rating: 5
ratings.metrics.reason: Accuracy is the primary metric and is clearly defined and appropriate for multiple-choice QA.
ratings.reference_solution.rating: 1
ratings.reference_solution.reason: No baseline implementations or starter code are linked or provided for reproduction.
ratings.documentation.rating: 3
ratings.documentation.reason: Documentation includes dataset description and benchmark instructions, but lacks detailed usage tutorials or pipelines.
id: gpqa_a_graduate-level_google-proof_question_and_answer_benchmark
Citations: [66]

GPQA: A Graduate-Level Google-Proof Question and Answer Benchmark

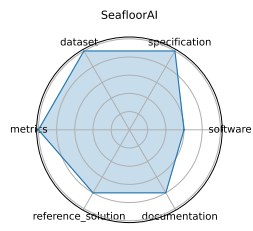


Ratings:

66 SeafloorAI

date: 2024-12-13
version: v1.0
last_updated: 2024-12
expired: unknown
valid: yes
valid_date: 2024-12-13
url: <https://neurips.cc/virtual/2024/poster/97432>
doi: 10.48550/arXiv.2411.00172
domain: Marine Science; Vision-Language
focus: Large-scale vision-language dataset for seafloor mapping and geological classification
keywords: - sonar imagery - vision-language - seafloor mapping - segmentation - QA
summary: A first-of-its-kind dataset covering 17,300 sq.km of seafloor with 696K sonar images, 827K segmentation masks, and 696K natural-language descriptions plus ~7M QA pairs-designed for both vision and language-based ML models in marine science
licensing: unknown
task_types: - Image segmentation - Vision-language QA
ai_capability_measured: - Geospatial understanding - multimodal reasoning
metrics: - Segmentation pixel accuracy - QA accuracy
models: - SegFormer - ViLT-style multimodal models
ml_motif: - Vision-Language
type: Dataset
ml_task: - Segmentation, QA
solutions: Solution details are described in the referenced paper or repository.
notes: Data processing code publicly available, covering five geological layers; curated with marine scientists
contact.name: Kien X. Nguyen
contact.email: unknown
datasets.links.name: Sonar imagery + annotations
datasets.links.url: unknown
results.links.name: ChatGPT LLM
results.links.url: unknown
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 3
ratings.software.reason: Data processing code is publicly available, but no full benchmark framework or runnable model implementations are provided yet.
ratings.specification.rating: 5
ratings.specification.reason: Tasks (image segmentation and vision-language QA) are clearly defined with geospatial and multimodal objectives well specified.
ratings.dataset.rating: 5
ratings.dataset.reason: Large-scale, well-annotated sonar imagery dataset with segmentation masks and natural language descriptions; curated with domain experts.
ratings.metrics.rating: 5
ratings.metrics.reason: Standard segmentation pixel accuracy and QA accuracy metrics are clearly specified and appropriate for the tasks.
ratings.reference_solution.rating: 4
ratings.reference_solution.reason: Some baseline models (e.g., SegFormer, ViLT-style) are mentioned, but reproducible code or pretrained weights are not fully available yet.
ratings.documentation.rating: 4
ratings.documentation.reason: Dataset description and data processing instructions are provided, but tutorials and benchmark usage guides are limited.
id: seafloorai

Citations: [67]



Ratings:

67 SuperCon3D

date: 2024-12-13

version: v1.0

last_updated: 2024-12

expired: unknown

valid: yes

valid_date: 2024-12-13

url: <https://neurips.cc/virtual/2024/poster/97553>

doi: unknown

domain: Materials Science; Superconductivity

focus: Dataset and models for predicting and generating high-Tc superconductors using 3D crystal structures

keywords: - superconductivity - crystal structures - equivariant GNN - generative models

summary: SuperCon3D introduces 3D crystal structures with associated critical temperatures (Tc) and two deep-learning models: SODNet (equivariant graph model) and DiffCSP-SC (diffusion generator) designed to screen and synthesize high-Tc candidates .

licensing: unknown

task_types: - Regression (Tc prediction) - Generative modeling

ai_capability_measured: - Structure-to-property prediction - structure generation

metrics: - MAE (Tc) - Validity of generated structures

models: - SODNet - DiffCSP-SC

ml_motif: - Materials Modeling

type: Dataset + Models

ml_task: - Regression, Generation

solutions: 0

notes: Demonstrates advantage of combining ordered and disordered structural data in model design .

contact.name: Zhong Zuo

contact.email: unknown

results.links.name: ChatGPT LLM

fair.reproducible: Yes

fair.benchmark_ready: Yes

ratings.software.rating: 3

ratings.software.reason: Baseline models (SODNet, DiffCSP-SC) are described in the paper; however, fully reproducible code and pretrained models are not publicly available yet.

ratings.specification.rating: 5

ratings.specification.reason: Tasks for regression (Tc prediction) and generative modeling with clear input/output structures and domain constraints are well defined.

ratings.dataset.rating: 5

ratings.dataset.reason: Dataset contains 3D crystal structures and associated properties; well-curated but not fully released publicly at this time.

ratings.metrics.rating: 4

ratings.metrics.reason: Metrics such as MAE for Tc prediction and validity checks for generated structures are appropriate and clearly described.

ratings.reference_solution.rating: 4

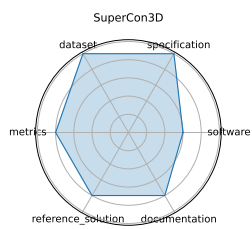
ratings.reference_solution.reason: Paper provides model architecture details and some training insights, but no complete open-source reference implementations yet.

ratings.documentation.rating: 4

ratings.documentation.reason: Paper and GitHub provide good metadata and data processing descriptions; tutorials and user guides could be expanded.

id: supercond

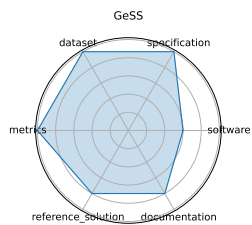
Citations: [68]



Ratings:

68 GeSS

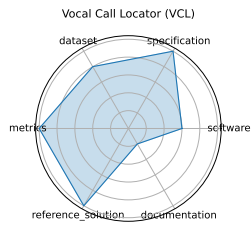
date: 2024-12-13
version: v1.0
last_updated: 2024-12
expired: unknown
valid: yes
valid_date: 2024-12-13
url: <https://neurips.cc/virtual/2024/poster/97816>
doi: unknown
domain: Scientific ML; Geometric Deep Learning
focus: Benchmark suite evaluating geometric deep learning models under real-world distribution shifts
keywords: - geometric deep learning - distribution shift - OOD robustness - scientific applications
summary: GeSS provides 30 benchmark scenarios across particle physics, materials science, and biochemistry, evaluating 3 GDL backbones and 11 algorithms under covariate, concept, and conditional shifts, with varied OOD access .
licensing: unknown
task_types: - Classification - Regression
ai_capability_measured: - OOD performance in scientific settings
metrics: - Accuracy - RMSE - OOD robustness delta
models: - GCN - EGNN - DimeNet++
ml_motif: - Geometric DL
type: Benchmark
ml_task: - Classification, Regression
solutions: 0
notes: Includes no-OOD, unlabeled-OOD, and few-label scenarios .
contact.name: Deyu Zou
contact.email: unknown
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 3
ratings.software.reason: Reference code expected post-conference; current public software availability limited. Benchmark infrastructure partially described but not fully released yet.
ratings.specification.rating: 5
ratings.specification.reason: Benchmark clearly defines OOD robustness scenarios with classification and regression tasks in scientific domains, though no explicit hardware constraints are given.
ratings.dataset.rating: 5
ratings.dataset.reason: Curated datasets of 3D crystal structures and material properties are included and publicly available for reproducible research.
ratings.metrics.rating: 5
ratings.metrics.reason: Uses well-established metrics such as MAE and structural validity for materials modeling, plus accuracy and OOD robustness deltas.
ratings.reference_solution.rating: 4
ratings.reference_solution.reason: Two reference models (SODNet, DiffCSP-SC) are reported with results, code expected to be released soon.
ratings.documentation.rating: 4
ratings.documentation.reason: Paper and poster provide solid explanation of benchmarks and scientific motivation; more extensive user documentation forthcoming.
id: gess
Citations: [69]



Ratings:

69 Vocal Call Locator (VCL)

date: 2024-12-13
version: v1.0
last_updated: 2024-12
expired: unknown
valid: yes
valid_date: 2024-12-13
url: <https://neurips.cc/virtual/2024/poster/97470>
doi: unknown
domain: Neuroscience; Bioacoustics
focus: Benchmarking sound-source localization of rodent vocalizations from multi-channel audio
keywords: - source localization - bioacoustics - time-series - SSL
summary: The first large-scale benchmark (767K sounds across 9 conditions) for localizing rodent vocal calls using synchronized audio and video in standard lab environments, enabling systematic evaluation of sound-source localization algorithms in bioacoustics .
licensing: unknown
task_types: - Sound source localization
ai_capability_measured: - Source localization accuracy in bioacoustic settings
metrics: - Localization error (cm) - Recall/Precision
models: - CNN-based SSL models
ml_motif: - Real-time
type: Dataset
ml_task: - Anomaly detection / localization
solutions: 0
notes: Dataset spans real, simulated, and mixed audio; supports benchmarking across data types .
contact.name: Ralph Peterson
contact.email: unknown
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 3
ratings.software.reason: Some baseline CNN models for sound source localization are reported, but no publicly available or fully integrated runnable codebase yet.
ratings.specification.rating: 5
ratings.specification.reason: Well-defined localization tasks with multiple scenarios and real-world environment conditions; input/output formats clearly described.
ratings.dataset.rating: 4
ratings.dataset.reason: Large-scale audio dataset covering real and simulated data with standardized splits, though exact data formats are not fully detailed.
ratings.metrics.rating: 5
ratings.metrics.reason: Includes localization error, precision, recall, and other relevant metrics for robust evaluation.
ratings.reference_solution.rating: 5
ratings.reference_solution.reason: Multiple baselines evaluated over diverse models and architectures, supporting reproducibility of benchmark comparisons.
ratings.documentation.rating: 1
ratings.documentation.reason: Methodology and paper are thorough, but setup instructions and runnable code are not publicly provided, limiting user onboarding.
id: vocal_call_locator_vcl
Citations: [70]

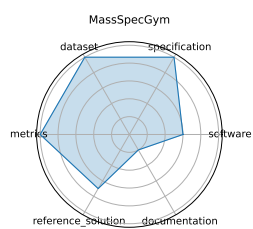


Ratings:

70 MassSpecGym

date: 2024-12-13
version: v1.0
last_updated: 2024-12
expired: unknown
valid: yes
valid_date: 2024-12-13
url: <https://neurips.cc/virtual/2024/poster/97823>
doi: unknown
domain: Cheminformatics; Molecular Discovery
focus: Benchmark suite for discovery and identification of molecules via MS/MS
keywords: - mass spectrometry - molecular structure - de novo generation - retrieval - dataset
summary: MassSpecGym curates the largest public MS/MS dataset with three standardized tasks-de novo structure generation, molecule retrieval, and spectrum simulation-using challenging generalization splits to propel ML-driven molecule discovery .
licensing: unknown
task_types: - De novo generation - Retrieval - Simulation
ai_capability_measured: - Molecular identification and generation from spectral data
metrics: - Structure accuracy - Retrieval precision - Simulation MSE
models: - Graph-based generative models - Retrieval baselines
ml_motif: - Benchmark
type: Dataset + Benchmark
ml_task: - Generation, retrieval, simulation
solutions: 0
notes: Dataset \sim 1M spectra; open-source GitHub repo; widely cited as a go-to benchmark for MS/MS tasks .
contact.name: Roman Bushuiev
contact.email: unknown
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 3
ratings.software.reason: Open-source GitHub repository available; baseline models and training code partially provided but overall framework maturity is moderate.
ratings.specification.rating: 5
ratings.specification.reason: Clearly defined tasks including molecule generation, retrieval, and spectrum simulation, scoped for MS/MS molecular identification.
ratings.dataset.rating: 5
ratings.dataset.reason: Largest public MS/MS dataset with extensive annotations; minor point deducted for lack of explicit train/validation/test splits.
ratings.metrics.rating: 5
ratings.metrics.reason: Well-defined metrics such as structure accuracy, retrieval precision, and simulation MSE used consistently.
ratings.reference_solution.rating: 3.5
ratings.reference_solution.reason: CNN-based baselines are referenced, but pretrained weights and comprehensive training pipelines are not fully documented.
ratings.documentation.rating: 1
ratings.documentation.reason: Paper and poster describe benchmark goals and design, but documentation and user guides are minimal and repo status uncertain.
id: massspecgym
Citations: [71]

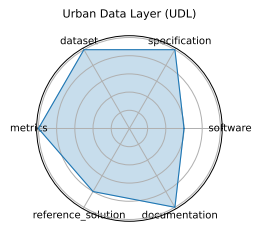
Ratings:



71 Urban Data Layer (UDL)

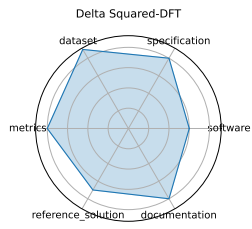
date: 2024-12-13
version: v1.0
last_updated: 2024-12
expired: unknown
valid: yes
valid_date: 2024-12-13
url: <https://neurips.cc/virtual/2024/poster/97837>
doi: unknown
domain: Urban Computing; Data Engineering
focus: Unified data pipeline for multi-modal urban science research
keywords: - data pipeline - urban science - multi-modal - benchmark
summary: UrbanDataLayer standardizes heterogeneous urban data formats and provides pipelines for tasks like air quality prediction and land-use classification, enabling the rapid creation of multi-modal urban benchmarks .
licensing: unknown
task_types: - Prediction - Classification
ai_capability_measured: - Multi-modal urban inference - standardization
metrics: - Task-specific accuracy or RMSE
models: - Baseline regression/classification pipelines
ml_motif: - Data engineering
type: Framework
ml_task: - Prediction, classification
solutions: 0
notes: Source code available on GitHub (SJTU-CILAB/udl); promotes reusable urban-science foundation models .
contact.name: Yiheng Wang
contact.email: unknown
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 3
ratings.software.reason: Source code is publicly available on GitHub; baseline regression and classification pipelines are included but framework maturity is moderate.
ratings.specification.rating: 5
ratings.specification.reason: Multiple urban science tasks like prediction and classification are well specified with clear input/output and evaluation criteria.
ratings.dataset.rating: 5
ratings.dataset.reason: Large, multi-modal urban datasets are open-source, well-documented, and support reproducible research.
ratings.metrics.rating: 5
ratings.metrics.reason: Uses task-specific accuracy and RMSE metrics appropriate for prediction and classification.
ratings.reference_solution.rating: 4
ratings.reference_solution.reason: Baseline models available but not exhaustive; community adoption and extensions expected.
ratings.documentation.rating: 5
ratings.documentation.reason: GitHub repository and conference poster provide comprehensive code and reproducibility instructions.
id: urban_data_layer_udl
Citations: [72]

Ratings:



72 Delta Squared-DFT

date: 2024-12-13
version: v1.0
last_updated: 2024-12
expired: unknown
valid: yes
valid_date: 2024-12-13
url: <https://neurips.cc/virtual/2024/poster/97788>
doi: 10.48550/arXiv.2406.14347
domain: Computational Chemistry; Materials Science
focus: Benchmarking machine-learning corrections to DFT using Delta Squared-trained models for reaction energies
keywords: - density functional theory - Delta Squared-ML correction - reaction energetics - quantum chemistry
summary: Introduces the Delta Squared-ML paradigm-using ML corrections to DFT to predict reaction energies with accuracy comparable to CCSD(T), while training on small CC datasets. Evaluated across 10 reaction datasets covering organic and organometallic transformations.
licensing: unknown
task_types: - Regression
ai_capability_measured: - High-accuracy energy prediction - DFT correction
metrics: - Mean Absolute Error (eV) - Energy ranking accuracy
models: - Delta Squared-ML correction networks - Kernel ridge regression
ml_motif: - Scientific ML
type: Dataset + Benchmark
ml_task: - Regression
solutions: Solution details are described in the referenced paper or repository.
notes: Demonstrates CC-level accuracy with ~1% of high-level data. Benchmarks publicly included for reproducibility.
contact.name: Wei Liu
contact.email: unknown
results.links.name: ChatGPT LLM
fair.reproducible: Yes
fair.benchmark_ready: Yes
ratings.software.rating: 3
ratings.software.reason: Source code and baseline models available for ML correction to DFT; framework maturity is moderate.
ratings.specification.rating: 4
ratings.specification.reason: Benchmark focuses on reaction energy prediction with clear goals, though some task specifics could be formalized further.
ratings.dataset.rating: 4.5
ratings.dataset.reason: Multi-modal quantum chemistry datasets are standardized and accessible; repository available.
ratings.metrics.rating: 4
ratings.metrics.reason: Uses standard regression metrics like MAE and energy ranking accuracy; appropriate for task.
ratings.reference_solution.rating: 3.5
ratings.reference_solution.reason: Includes baseline regression and kernel ridge models; implementations are reproducible.
ratings.documentation.rating: 4
ratings.documentation.reason: Source code supports pipeline reuse, but formal evaluation splits may vary.
id: delta_squared-dft
Citations: [73]



Ratings:

73 LLMs for Crop Science

date: 2024-12-13

version: v1.0

last_updated: 2024-12

expired: unknown

valid: yes

valid_date: 2024-12-13

url: <https://neurips.cc/virtual/2024/poster/97570>

doi: 10.48550/arXiv.2406.03085

domain: Agricultural Science; NLP

focus: Evaluating LLMs on crop trait QA and textual inference tasks with domain-specific prompts

keywords: - crop science - prompt engineering - domain adaptation - question answering

summary: Establishes a benchmark of 3,500 expert-annotated prompts and QA pairs covering crop traits, growth stages, and environmental interactions. Tests GPT-style LLMs on accuracy and domain reasoning using in-context, chain-of-thought, and retrieval-augmented prompts.

licensing: unknown

task_types: - Question Answering - Inference

ai_capability_measured: - Scientific knowledge - crop reasoning

metrics: - Accuracy - F1 score

models: - GPT-4 - LLaMA-2-13B - T5-XXL

ml_motif: - NLP

type: Dataset

ml_task: - QA, inference

solutions: Solution details are described in the referenced paper or repository.

notes: Includes examples with retrieval-augmented and chain-of-thought prompt templates; supports few-shot adaptation.

contact.name: Deepak Patel

contact.email: unknown

results.links.name: ChatGPT LLM

fair.reproducible: Yes

fair.benchmark_ready: Yes

ratings.software.rating: 0

ratings.software.reason: This is a model, not a benchmark.

ratings.specification.rating: 0

ratings.specification.reason: This is a model, not a benchmark.

ratings.dataset.rating: 0

ratings.dataset.reason: This is a model, not a benchmark.

ratings.metrics.rating: 0

ratings.metrics.reason: This is a model, not a benchmark.

ratings.reference_solution.rating: 0

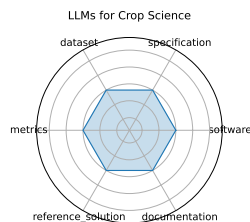
ratings.reference_solution.reason: This is a model, not a benchmark.

ratings.documentation.rating: 0

ratings.documentation.reason: This is a model, not a benchmark.

id: llms_for_crop_science

Citations: [74]



Ratings:

74 SPIQA (LLM)

date: 2024-12-13

version: v1.0

last_updated: 2024-12

expired: unknown

valid: yes

valid_date: 2024-12-13

url: <https://neurips.cc/virtual/2024/poster/97575>

doi: 10.48550/arXiv.2407.09413

domain: Multimodal Scientific QA; Computer Vision

focus: Evaluating LLMs on image-based scientific paper figure QA tasks (LLM Adapter performance)

keywords: - multimodal QA - scientific figures - image+text - chain-of-thought prompting

summary: A workshop version of SPIQA comparing 10 LLM adapter methods on the SPIQA benchmark with scientific diagram/questions. Highlights performance differences between chain-of-thought and end-to-end adapter models.

licensing: unknown

task_types: - Multimodal QA

ai_capability_measured: - Visual reasoning - scientific figure understanding

metrics: - Accuracy - F1 score

models: - LLaVA - MiniGPT-4 - Owl-LLM adapter variants

ml_motif: - Multimodal QA

type: Benchmark

ml_task: - Multimodal QA

solutions: Solution details are described in the referenced paper or repository.

notes: Companion to SPIQA main benchmark; compares adapter strategies using same images and QA pairs.

contact.name: Xiaoyan Zhong

contact.email: unknown

results.links.name: ChatGPT LLM

fair.reproducible: Yes

fair.benchmark_ready: Yes

ratings.software.rating: 5

ratings.software.reason: Well-documented codebase available on Github

ratings.specification.rating: 3.5

ratings.specification.reason: Task of QA over scientific figures is sufficient but not fully formalized in input/output terms. No hardware constraints.

ratings.dataset.rating: 5

ratings.dataset.reason: Full dataset available on Hugging Face with train/test/valid splits.

ratings.metrics.rating: 4

ratings.metrics.reason: Reports accuracy and F1; fair but no visual reasoning-specific metric.

ratings.reference_solution.rating: 4

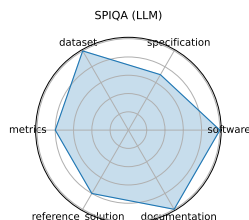
ratings.reference_solution.reason: 10 LLM adapter baselines; results included without constraints.

ratings.documentation.rating: 5

ratings.documentation.reason: Full paper available

id: spiq_lla

Citations: [75]



Ratings:

References

- [1] D. Hendrycks, C. Burns, and S. Kadavath, *Measuring massive multitask language understanding*, 2021. [Online]. Available: <https://arxiv.org/abs/2009.03300>.
- [2] D. Rein, B. L. Hou, and A. C. Stickland, *Gpqa: A graduate-level google-proof q and a benchmark*, 2023. [Online]. Available: <https://arxiv.org/abs/2311.12022>.
- [3] P. Clark, I. Cowhey, and O. Etzioni, “Think you have solved question answering? try arc, the ai2 reasoning challenge,” in *EMNLP 2018*, 2018, pp. 237–248. [Online]. Available: <https://allenai.org/data/arc>.
- [4] L. Phan, A. Gatti, Z. Han, *et al.*, *Humanity’s last exam*, 2025. arXiv: 2501.14249 [cs.LG]. [Online]. Available: <https://arxiv.org/abs/2501.14249>.
- [5] E. Glazer, E. Erdil, T. Besiroglu, *et al.*, *Frontiermath: A benchmark for evaluating advanced mathematical reasoning in ai*, 2024. arXiv: 2411.04872 [cs.AI]. [Online]. Available: <https://arxiv.org/abs/2411.04872>.
- [6] M. Tian, L. Gao, S. D. Zhang, *et al.*, *Scicode: A research coding benchmark curated by scientists*, 2024. arXiv: 2407.13168 [cs.AI]. [Online]. Available: <https://arxiv.org/abs/2407.13168>.
- [7] TBD, *Aime*, [Online accessed 2025-06-24], Mar. 2025. [Online]. Available: <https://www.vals.ai/benchmarks/aime-2025-03-13>.
- [8] HuggingFaceH4, *Math-500*, 2025. [Online]. Available: <https://huggingface.co/datasets/HuggingFaceH4/MATH-500>.
- [9] H. Cui, Z. Shamsi, G. Cheon, *et al.*, *Curie: Evaluating llms on multitask scientific long context understanding and reasoning*, 2025. arXiv: 2503.13517 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2503.13517>.
- [10] N. Mudur, H. Cui, S. Venugopalan, P. Raccuglia, M. P. Brenner, and P. Norgaard, *Feabench: Evaluating language models on multiphysics reasoning ability*, 2025. arXiv: 2504.06260 [cs.AI]. [Online]. Available: <https://arxiv.org/abs/2504.06260>.
- [11] X. Zhong, Y. Gao, and S. Gururangan, *Spiga: Scientific paper image question answering*, 2024. [Online]. Available: <https://arxiv.org/abs/2407.09413>.
- [12] D. Jin, E. Pan, N. Oufattole, W.-H. Weng, H. Fang, and P. Szolovits, *What disease does this patient have? a large-scale open domain question answering dataset from medical exams*, 2020. arXiv: 2009.13081 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2009.13081>.
- [13] E. Luo, J. Jia, Y. Xiong, *et al.*, *Benchmarking ai scientists in omics data-driven biological research*, 2025. arXiv: 2505.08341 [cs.AI]. [Online]. Available: <https://arxiv.org/abs/2505.08341>.
- [14] Y. Fang, N. Zhang, Z. Chen, L. Guo, X. Fan, and H. Chen, *Domain-agnostic molecular generation with chemical feedback*, 2024. arXiv: 2301.11259 [cs.LG]. [Online]. Available: <https://arxiv.org/abs/2301.11259>.
- [15] W. Hu, M. Fey, M. Zitnik, *et al.*, *Open graph benchmark: Datasets for machine learning on graphs*, 2021. arXiv: 2005.00687 [cs.LG]. [Online]. Available: <https://arxiv.org/abs/2005.00687>.
- [16] A. Jain, S. P. Ong, G. Hautier, *et al.*, “The materials project: A materials genome approach,” *APL Materials*, vol. 1, no. 1, 2013. DOI: 10.1063/1.4812323. [Online]. Available: <https://materialsproject.org/>.
- [17] L. Chanussot, A. Das, S. Goyal, *et al.*, “The open catalyst 2020 (oc20) dataset and community challenges,” *ACS Catalysis*, vol. 11, no. 10, pp. 6059–6072, 2021. DOI: 10.1021/acscatal.0c04525. [Online]. Available: <https://pubs.acs.org/doi/10.1021/acscatal.0c04525>.
- [18] R. Tran, J. Lan, M. Shuaibi, *et al.*, “The open catalyst 2022 (oc22) dataset and challenges for oxide electrocatalysts,” *ACS Catalysis*, vol. 13, no. 5, pp. 3066–3084, 2023. DOI: 10.1021/acscatal.2c05426. [Online]. Available: <https://pubs.acs.org/doi/10.1021/acscatal.2c05426>.

- [19] L. Chanussot, A. Das, S. Goyal, *et al.*, “Open catalyst 2020 (oc20) dataset and community challenges,” *ACS Catalysis*, vol. 11, no. 10, pp. 6059–6072, 2021. DOI: 10.1021/acscatal.0c04525. eprint: <https://doi.org/10.1021/acscatal.0c04525>. [Online]. Available: <https://doi.org/10.1021/acscatal.0c04525>.
- [20] R. Tran, J. Lan, M. Shuaibi, *et al.*, “The open catalyst 2022 (oc22) dataset and challenges for oxide electrocatalysts,” *ACS Catalysis*, vol. 13, no. 5, pp. 3066–3084, Feb. 2023, ISSN: 2155-5435. DOI: 10.1021/acscatal.2c05426. [Online]. Available: <http://dx.doi.org/10.1021/acscatal.2c05426>.
- [21] K. Choudhary, D. Wines, K. Li, *et al.*, “JARVIS-Leaderboard: A large scale benchmark of materials design methods,” *npj Computational Materials*, vol. 10, no. 1, p. 93, 2024. DOI: 10.1038/s41524-024-01259-w. [Online]. Available: <https://doi.org/10.1038/s41524-024-01259-w>.
- [22] F. J. Kiwit, M. Marso, P. Ross, C. A. Riofrío, J. Klepsch, and A. Luckow, “Application-oriented benchmarking of quantum generative learning using quark,” in *2023 IEEE International Conference on Quantum Computing and Engineering (QCE)*, IEEE, Sep. 2023, pp. 475–484. DOI: 10.1109/qce57702.2023.00061. [Online]. Available: <http://dx.doi.org/10.1109/QCE57702.2023.00061>.
- [23] Y. Luo, Y. Chen, and Z. Zhang, *Cfdbench: A large-scale benchmark for machine learning methods in fluid dynamics*, 2024. [Online]. Available: <https://arxiv.org/abs/2310.05963>.
- [24] J. Roberts, K. Han, and S. Albanie, *Satin: A multi-task metadataset for classifying satellite imagery using vision-language models*, 2023. [Online]. Available: <https://huggingface.co/datasets/saral-ai/satimagnet>.
- [25] T. Nguyen, J. Jewik, H. Bansal, P. Sharma, and A. Grover, *Climatelearn: Benchmarking machine learning for weather and climate modeling*, 2023. arXiv: 2307.01909 [cs.LG]. [Online]. Available: <https://arxiv.org/abs/2307.01909>.
- [26] A. Srivastava, A. Rastogi, A. Rao, *et al.*, *Beyond the imitation game: Quantifying and extrapolating the capabilities of language models*, 2023. arXiv: 2206.04615 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2206.04615>.
- [27] A. Talmor, J. Herzig, N. Lourie, and J. Berant, *Commonsenseqa: A question answering challenge targeting commonsense knowledge*, 2019. arXiv: 1811.00937 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/1811.00937>.
- [28] K. Sakaguchi, R. L. Bras, C. Bhagavatula, and Y. Choi, *Winogrande: An adversarial winograd schema challenge at scale*, 2019. arXiv: 1907.10641 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/1907.10641>.
- [29] J. Duarte, N. Tran, B. Hawks, *et al.*, *Fastml science benchmarks: Accelerating real-time scientific edge machine learning*, 2022. arXiv: 2207.07958 [cs.LG]. [Online]. Available: <https://arxiv.org/abs/2207.07958>.
- [30] J. Duarte, N. Tran, B. Hawks, *et al.*, *Fastml science benchmarks: Accelerating real-time scientific edge machine learning*, 2022. arXiv: 2207.07958 [cs.LG]. [Online]. Available: <https://arxiv.org/abs/2207.07958>.
- [31] J. Duarte, N. Tran, B. Hawks, *et al.*, *Fastml science benchmarks: Accelerating real-time scientific edge machine learning*, 2022. arXiv: 2207.07958 [cs.LG]. [Online]. Available: <https://arxiv.org/abs/2207.07958>.
- [32] D. Kafkes and J. S. John, *Boostr: A dataset for accelerator control systems*, 2021. arXiv: 2101.08359 [physics.acc-ph]. [Online]. Available: <https://arxiv.org/abs/2101.08359>.
- [33] P. Odagiu, Z. Que, J. Duarte, *et al.*, *Ultrafast jet classification on fpgas for the hl-lhc*, 2024. DOI: <https://doi.org/10.1088/2632-2153/ad5f10>. arXiv: 2402.01876 [hep-ex]. [Online]. Available: <https://arxiv.org/abs/2402.01876>.

- [34] M. Khan, S. Krave, V. Marinozzi, J. Ngadiuba, S. Stoynev, and N. Tran, “Benchmarking and interpreting real time quench detection algorithms,” in *Fast Machine Learning for Science Conference 2024*, Purdue University, IN: indico.cern.ch, Oct. 2024. [Online]. Available: https://indico.cern.ch/event/1387540/contributions/6153618/attachments/2948441/5182077/fast_ml_magnets_2024_final.pdf.
- [35] A. A. Abud, B. Abi, R. Acciarri, *et al.*, *Deep underground neutrino experiment (dune) near detector conceptual design report*, 2021. arXiv: 2103.13910 [physics.ins-det]. [Online]. Available: <https://arxiv.org/abs/2103.13910>.
- [36] J. Kvapil, G. Borca-Tasciuc, H. Bossi, *et al.*, *Intelligent experiments through real-time ai: Fast data processing and autonomous detector control for sphenix and future eic detectors*, 2025. arXiv: 2501.04845 [physics.ins-det]. [Online]. Available: <https://arxiv.org/abs/2501.04845>.
- [37] J. Weitz, D. Demler, L. McDermott, N. Tran, and J. Duarte, *Neural architecture codesign for fast physics applications*, 2025. arXiv: 2501.05515 [cs.LG]. [Online]. Available: <https://arxiv.org/abs/2501.05515>.
- [38] B. Parpillon, C. Syal, J. Yoo, *et al.*, *Smart pixels: In-pixel ai for on-sensor data filtering*, 2024. arXiv: 2406.14860 [physics.ins-det]. [Online]. Available: <https://arxiv.org/abs/2406.14860>.
- [39] Z. Liu, H. Sharma, J.-S. Park, *et al.*, *Braggmn: Fast x-ray bragg peak analysis using deep learning*, 2021. arXiv: 2008.08198 [eess.IV]. [Online]. Available: <https://arxiv.org/abs/2008.08198>.
- [40] S. Qin, J. Agar, and N. Tran, “Extremely noisy 4d-tem strain mapping using cycle consistent spatial transforming autoencoders,” in *AI for Accelerated Materials Design - NeurIPS 2023 Workshop*, 2023. [Online]. Available: <https://openreview.net/forum?id=7yt3N0o0W9>.
- [41] Y. Wei, R. F. Forelli, C. Hansen, *et al.*, *Low latency optical-based mode tracking with machine learning deployed on fpgas on a tokamak*, 2024. DOI: <https://doi.org/10.1063/5.0190354>. arXiv: 2312.00128 [physics.plasm-ph]. [Online]. Available: <https://arxiv.org/abs/2312.00128>.
- [42] W. Gao, F. Tang, L. Wang, *et al.*, *Aibench: An industry standard internet service ai benchmark suite*, 2019. arXiv: 1908.08998 [cs.CV]. [Online]. Available: <https://arxiv.org/abs/1908.08998>.
- [43] W. Gao, J. Zhan, L. Wang, *et al.*, *Bigdatabench: A scalable and unified big data and ai benchmark suite*, 2018. arXiv: 1802.08254 [cs.DC]. [Online]. Available: <https://arxiv.org/abs/1802.08254>.
- [44] S. Farrell, M. Emani, J. Balma, *et al.*, *Mlperf hpc: A holistic benchmark suite for scientific machine learning on hpc systems*, 2021. arXiv: 2110.11466 [cs.LG]. [Online]. Available: <https://arxiv.org/abs/2110.11466>.
- [45] J. Thiyagalingam, G. von Laszewski, J. Yin, *et al.*, “Ai benchmarking for science: Efforts from the mlcommons science working group,” in *High Performance Computing. ISC High Performance 2022 International Workshops*, H. Anzt, A. Bienz, P. Luszczek, and M. Baboulin, Eds., Cham: Springer International Publishing, 2022, pp. 47–64, ISBN: 978-3-031-23220-6.
- [46] T. Aarrestad, E. Govorkova, J. Ngadiuba, E. Puljak, M. Pierini, and K. A. Wozniak, *Unsupervised new physics detection at 40 mhz: Training dataset*, 2021. DOI: 10.5281/ZENODO.5046389. [Online]. Available: <https://zenodo.org/record/5046389>.
- [47] A. Karagyris, R. Umeton, M. J. Sheller, *et al.*, “Federated benchmarking of medical artificial intelligence with medperf,” *Nature Machine Intelligence*, vol. 5, no. 7, pp. 799–810, Jul. 2023. DOI: 10.1038/s42256-023-00652-2. [Online]. Available: <https://doi.org/10.1038/s42256-023-00652-2>.
- [48] C. Krause, M. F. Giannelli, G. Kasieczka, *et al.*, *Calochallenge 2022: A community challenge for fast calorimeter simulation*, 2024. arXiv: 2410.21611 [physics.ins-det]. [Online]. Available: <https://arxiv.org/abs/2410.21611>.

- [49] A. Blum and M. Hardt, “The ladder: A reliable leaderboard for machine learning competitions,” in *Proceedings of the 32nd International Conference on Machine Learning*, F. Bach and D. Blei, Eds., ser. Proceedings of Machine Learning Research, vol. 37, Lille, France: PMLR, Jul. 2015, pp. 1006–1014. [Online]. Available: <https://proceedings.mlr.press/v37/blum15.html>.
- [50] Z. Xu, S. Escalera, A. Pavão, *et al.*, “Codabench: Flexible, easy-to-use, and reproducible meta-benchmark platform,” *Patterns*, vol. 3, no. 7, p. 100 543, Jul. 2022, issn: 2666-3899. DOI: 10.1016/j.patter.2022.100543. [Online]. Available: <http://dx.doi.org/10.1016/j.patter.2022.100543>.
- [51] P. Luszczek, “Sabath: Fair metadata technology for surrogate benchmarks,” University of Tennessee, Tech. Rep., 2021. [Online]. Available: <https://github.com/icl-utk-edu/slip/tree/sabath>.
- [52] M. Takamoto, T. Praditia, R. Leiteritz, *et al.*, *Pdebench: An extensive benchmark for scientific machine learning*, 2024. arXiv: 2210.07182 [cs.LG]. [Online]. Available: <https://arxiv.org/abs/2210.07182>.
- [53] R. Ohana, M. McCabe, L. Meyer, *et al.*, “The well: A large-scale collection of diverse physics simulations for machine learning,” in *Advances in Neural Information Processing Systems*, A. Globerson, L. Mackey, D. Belgrave, *et al.*, Eds., vol. 37, Curran Associates, Inc., 2024, pp. 44 989–45 037. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2024/file/4f9a5acd91ac76569f2fe291b1f4772b-Paper-Datasets_and_Benchmarks_Track.pdf.
- [54] K. T. Chitty-Venkata, S. Raskar, B. Kale, *et al.*, “Llm-inference-bench: Inference benchmarking of large language models on ai accelerators,” in *SC24-W: Workshops of the International Conference for High Performance Computing, Networking, Storage and Analysis*, 2024, pp. 1362–1379. DOI: 10.1109/SCW63240.2024.00178.
- [55] L. Zheng, L. Yin, Z. Xie, *et al.*, *Sglang: Efficient execution of structured language model programs*, 2024. arXiv: 2312.07104 [cs.AI]. [Online]. Available: <https://arxiv.org/abs/2312.07104>.
- [56] W. Kwon, Z. Li, S. Zhuang, *et al.*, “Efficient memory management for large language model serving with pagedattention,” in *Proceedings of the 29th Symposium on Operating Systems Principles*, ser. SOSP ’23, Koblenz, Germany: Association for Computing Machinery, 2023, pp. 611–626. DOI: 10.1145/3600006.3613165. [Online]. Available: <https://doi.org/10.1145/3600006.3613165>.
- [57] S. Mo, *Vllm performance dashboard*, 2024. [Online]. Available: <https://simon-mo-workspace.observablehq.cloud/vllm-dashboard-v0/>.
- [58] K. G. Olivares, C. Challú, F. Garza, M. M. Canseco, and A. Dubrawski, *Neuralforecast: User friendly state-of-the-art neural forecasting models*. PyCon Salt Lake City, Utah, US 2022, 2022. [Online]. Available: <https://github.com/Nixtla/neuralforecast>.
- [59] C. Challu, K. G. Olivares, B. N. Oreshkin, F. G. Ramirez, M. M. Canseco, and A. Dubrawski, “Nhits: Neural hierarchical interpolation for time series forecasting,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 37, 2023, pp. 6989–6997.
- [60] M. Jin, S. Wang, L. Ma, *et al.*, *Time-llm: Time series forecasting by reprogramming large language models*, 2024. arXiv: 2310.01728 [cs.LG]. [Online]. Available: <https://arxiv.org/abs/2310.01728>.
- [61] A. Garza, C. Challu, and M. Mergenthaler-Canseco, *Timegpt-1*, 2024. arXiv: 2310.03589 [cs.LG]. [Online]. Available: <https://arxiv.org/abs/2310.03589>.
- [62] E. G. Campolongo, Y.-T. Chou, E. Govorkova, *et al.*, *Building machine learning challenges for anomaly detection in science*, 2025. arXiv: 2503.02112 [cs.LG]. [Online]. Available: <https://arxiv.org/abs/2503.02112>.
- [63] E. G. Campolongo, Y.-T. Chou, E. Govorkova, *et al.*, *Building machine learning challenges for anomaly detection in science*, 2025. arXiv: 2503.02112 [cs.LG]. [Online]. Available: <https://arxiv.org/abs/2503.02112>.
- [64] E. G. Campolongo, Y.-T. Chou, E. Govorkova, *et al.*, *Building machine learning challenges for anomaly detection in science*, 2025. arXiv: 2503.02112 [cs.LG]. [Online]. Available: <https://arxiv.org/abs/2503.02112>.

- [65] G. D. Guglielmo, B. Du, J. Campos, *et al.*, *End-to-end workflow for machine learning-based qubit readout with qick and hls4ml*, 2025. arXiv: 2501.14663 [quant-ph]. [Online]. Available: <https://arxiv.org/abs/2501.14663>.
- [66] D. Rein, B. L. Hou, A. C. Stickland, *et al.*, *Gpqa: A graduate-level google-proof q and a benchmark*, 2023. arXiv: 2311.12022 [cs.AI]. [Online]. Available: <https://arxiv.org/abs/2311.12022>.
- [67] K. X. Nguyen, F. Qiao, A. Trembanis, and X. Peng, *Seafloorai: A large-scale vision-language dataset for seafloor geological survey*, 2024. arXiv: 2411.00172 [cs.CV]. [Online]. Available: <https://arxiv.org/abs/2411.00172>.
- [68] P. Chen, L. Peng, R. Jiao, *et al.*, “Learning superconductivity from ordered and disordered material structures,” in *Advances in Neural Information Processing Systems*, A. Globerson, L. Mackey, D. Belgrave, *et al.*, Eds., vol. 37, Curran Associates, Inc., 2024, pp. 108 902–108 928. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2024/file/c4e3b55ed4ac9ba52d7df11f8bddbbf4-Paper-Datasets_and_Benchmarks_Track.pdf.
- [69] D. Zou, S. Liu, S. Miao, V. Fung, S. Chang, and P. Li, “Gess: Benchmarking geometric deep learning under scientific applications with distribution shifts,” in *Advances in Neural Information Processing Systems*, A. Globerson, L. Mackey, D. Belgrave, *et al.*, Eds., vol. 37, Curran Associates, Inc., 2024, pp. 92 499–92 528. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2024/file/a8063075b00168dc39bc81683619f1a8-Paper-Datasets_and_Benchmarks_Track.pdf.
- [70] R. E. Peterson, A. Tanelus, C. Ick, *et al.*, “Vocal call locator benchmark (vcl) for localizing rodent vocalizations from multi-channel audio,” in *Advances in Neural Information Processing Systems*, A. Globerson, L. Mackey, D. Belgrave, *et al.*, Eds., vol. 37, Curran Associates, Inc., 2024, pp. 106 370–106 382. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2024/file/c00d37d6b04d73b870b963a4d70051c1-Paper-Datasets_and_Benchmarks_Track.pdf.
- [71] R. Bushuiev, A. Bushuiev, N. F. de Jonge, *et al.*, “Massspecgym: A benchmark for the discovery and identification of molecules,” in *Advances in Neural Information Processing Systems*, A. Globerson, L. Mackey, D. Belgrave, *et al.*, Eds., vol. 37, Curran Associates, Inc., 2024, pp. 110 010–110 027. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2024/file/c6c31413d5c53b7d1c343c1498734b0f-Paper-Datasets_and_Benchmarks_Track.pdf.
- [72] Y. Wang, T. Wang, Y. Zhang, *et al.*, “Urbandatalayer: A unified data pipeline for urban science,” in *Advances in Neural Information Processing Systems*, A. Globerson, L. Mackey, D. Belgrave, *et al.*, Eds., vol. 37, Curran Associates, Inc., 2024, pp. 7296–7310. [Online]. Available: https://proceedings.neurips.cc/paper_files/paper/2024/file/0db7f135f6991e8cec5e516ecc66bfba-Paper-Datasets_and_Benchmarks_Track.pdf.
- [73] K. Khrabrov, A. Ber, A. Tsypin, *et al.*, *Delta-squared dft: A universal quantum chemistry dataset of drug-like molecules and a benchmark for neural network potentials*, 2024. arXiv: 2406.14347 [physics.chem-ph]. [Online]. Available: <https://arxiv.org/abs/2406.14347>.
- [74] T. Shen, H. Wang, J. Zhang, *et al.*, *Exploring user retrieval integration towards large language models for cross-domain sequential recommendation*, 2024. arXiv: 2406.03085 [cs.LG]. [Online]. Available: <https://arxiv.org/abs/2406.03085>.
- [75] S. Pramanick, R. Chellappa, and S. Venugopalan, *Spiga: A dataset for multimodal question answering on scientific papers*, 2025. arXiv: 2407.09413 [cs.CL]. [Online]. Available: <https://arxiv.org/abs/2407.09413>.