

# OMNIFLOW: A Physics-Grounded Multimodal Agent for Generalized Scientific Reasoning

Anonymous ACL submission

## Abstract

Large Language Models (LLMs) have demonstrated exceptional logical reasoning capabilities but frequently struggle with the continuous spatiotemporal dynamics governed by Partial Differential Equations (PDEs), often resulting in non-physical hallucinations. Existing approaches typically resort to costly, domain-specific fine-tuning, which severely limits cross-domain generalization and interpretability. To bridge this gap, we propose OMNIFLOW, a neuro-symbolic architecture designed to ground frozen multimodal LLMs in fundamental physical laws without requiring domain-specific parameter updates. OMNIFLOW introduces a novel *Semantic-Symbolic Alignment* mechanism that projects high-dimensional flow tensors into topological linguistic descriptors, enabling the model to perceive physical structures rather than raw pixel values. Furthermore, we construct a Physics-Guided Chain-of-Thought (PG-CoT) workflow that orchestrates reasoning through dynamic constraint injection (e.g., mass conservation) and iterative reflexive verification. We evaluate OMNIFLOW on a comprehensive benchmark spanning microscopic turbulence, theoretical Navier-Stokes equations, and macroscopic global weather forecasting. Empirical results demonstrate that OMNIFLOW significantly outperforms traditional deep learning baselines in zero-shot generalization and few-shot adaptation tasks. Crucially, it offers transparent, physically consistent reasoning reports, marking a paradigm shift from black-box fitting to interpretable scientific reasoning. Our code is available at <https://anonymous.4open.science/r/OMNIFLOW-061B>.

## 1 Introduction

Large Language Models (LLMs) (Chang et al., 2024; Naveed et al., 2023; Zhao et al., 2023) have demonstrated exceptional symbolic reasoning, code generation, and mathematical problem-solving capabilities. However, when applied to

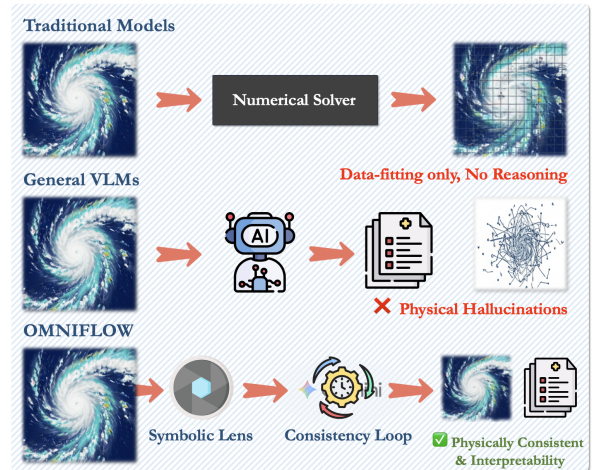


Figure 1: **Comparison of reasoning paradigms.** Traditional models (top) are non-interpretable black-boxes. General VLMs (middle) suffer from physical hallucinations. OMNIFLOW (bottom) integrates a *Symbolic Lens* and *Consistency Loop* to deliver grounded forecasts and expert reports, uniting numerical precision with logical reasoning.

the physical world governed by Partial Differential Equations (PDEs) (Chen and Shaw, 2001; Evans, 2010), particularly systems involving continuous spatiotemporal dynamics (Yu et al., 2018; Mohan et al., 2020) such as turbulence evolution (Davidson, 2015; Pope, 2001) or global weather forecasting (Rasp et al., 2020; Gao et al., 2022b; Wu et al., 2025b; Gao et al., 2025), LLMs often exhibit significant limitations. Lacking physical grounding, existing multimodal models struggle to comprehend the topological structures within high-dimensional fluid data, frequently resulting in non-physical hallucinations that violate fundamental physical common sense.

Historically, two main paradigms have addressed this challenge. ❶. The first involves specialized deep learning models (e.g., FNO (Li et al., 2020), GraphCast (Lam et al., 2023)), which serve as surrogates for numerical solvers. While accurate,

064 these models function primarily as data fitters, lack- 116  
065 ing cross-domain generalization capabilities and 117  
066 interpretability. ②. The second paradigm involves 118  
067 fine-tuning LLMs with large-scale scientific data. 119  
068 However, fine-tuning incurs high computational 120  
069 costs, often leads to catastrophic forgetting of gen- 121  
070 eral knowledge, and still fails to guarantee strict 122  
071 adherence to physical conservation laws, such as 123  
072 mass or momentum conservation (Du et al., 2024). 124  
073 As illustrated in Figure 1, even advanced Vision- 125  
074 Language Models (VLMs) fall into this second 126  
075 category’s trap: they interpret scientific imagery as 127  
076 semantic patterns rather than discrete solutions to 128  
077 PDEs, leading to visually plausible but physically 129  
078 invalid outputs.

079 We argue that the inability of LLMs to handle 130  
080 physical problems stems not from a lack of internal 131  
081 knowledge, but from a lack of modal alignment and 132  
082 logical constraint. Existing general-purpose VLMs 133  
083 fundamentally lack the numerical precision and 134  
084 physical inductive bias required for spatiotemporal 135  
085 fluid forecasting. Consequently, relying solely on 136  
086 standard VLMs leads to prognostic outputs that 137  
087 strictly violate conservation laws. To address this, 138  
088 we explore a novel paradigm: instead of modifying 139  
089 model parameters, we design a Cognitive Archite- 140  
090 cture to decouple physical computation from cogni- 141  
091 tive reasoning, aligning frozen LLMs with rigorous 142  
092 physical rules via neuro-symbolic collaboration. 143  
093

094 To this end, we propose OMNIFLOW, a physics- 144  
095 grounded agent framework for generalized fluid 145  
096 dynamics reasoning. As shown in Figure 1, OM- 146  
097 NIFLOW abandons the traditional black-box pre- 147  
098 diction mode (Raissi et al., 2019; Li et al., 2020; 148  
099 Bi et al., 2023; Wu et al., 2024) in favor of a trans- 149  
100 parent reasoning workflow. First, addressing the 150  
101 heterogeneity of multimodal inputs (Dosovitskiy 151  
102 et al., 2021; Liu et al., 2024), we design a Visual 152  
103 Symbolic Projector (the *Symbolic Lens* in Figure 1). 153  
104 This module translates raw flow fields (e.g., satel- 154  
105 lite typhoon imagery) into semantic tokens contain- 155  
106 ing vector field features and topological skeletons, 156  
107 achieving alignment between continuous data and 157  
108 discrete symbols. Second, we introduce a Physics- 158  
109 Guided Chain-of-Thought (PG-CoT) (Wei et al., 159  
110 2022). Within the reasoning engine, the agent oper- 160  
111 ates an In-Context Reflexive Loop (the *Consistency* 161  
112 *Loop* in Figure 1): it dynamically retrieves external 162  
113 physical knowledge (e.g., Navier-Stokes equations) 163  
114 and executes a Consistency Check during genera- 164  
115 tion (Lewis et al., 2020). Upon detecting a physical 165  
violation (e.g., a trajectory violating inertial con-

straints), a Critic module forces the model to roll 116  
back and self-correct. 117

We ground agentic reasoning in physical laws 118  
rather than linguistic plausibility, using symbolic 119  
verification to detect and rectify invalid reasoning 120  
steps. The main contributions of this paper are as 121  
follows: 122

1. **Architectural Innovation:** We propose the 123  
first VLLM training-free framework for general- 124  
ized fluid physical reasoning. Through a neuro- 125  
symbolic mechanism, OMNIFLOW successfully ac- 126  
tivates the reasoning potential of frozen LLMs for 127  
complex scientific computing tasks without costly 128  
parameter updates. 129

2. **Generalization:** We evaluate OMNIFLOW on 130  
three distinct physical benchmarks spanning micro- 131  
scopic turbulence, theoretical Navier-Stokes equa- 132  
tions, and macroscopic global weather forecast- 133  
ing. Experiments demonstrate that in a Zero-Shot 134  
setting, OMNIFLOW not only adapts to different 135  
governing equations but also achieves prediction 136  
accuracy comparable to specialized deep learning 137  
models. 138

3. **Interpretability:** Unlike the mute numerical 139  
outputs of traditional methods, OMNIFLOW gener- 140  
ates structured analysis reports containing physi- 141  
cal grounding, risk assessment, and decision logic, 142  
providing a new interaction paradigm for scientific 143  
discovery and decision support. 144

## 2 Related Work 145

**Deep Learning for Fluid Dynamics** Deep learn- 146  
ing has emerged as a powerful paradigm for ac- 147  
celerating fluid dynamics simulations (Kutz, 2017; 148  
Brunton et al., 2020), traditionally reliant on com- 149  
putationally expensive numerical solvers. Early 150  
data-driven approaches utilized Convolutional Neu- 151  
ral Networks (CNNs) to approximate flow fields on 152  
fixed grids (Shi et al., 2015; He et al., 2016; Raonic 153  
et al., 2023). More recently, Physics-Informed Neu- 154  
ral Networks (PINNs) (Raissi et al., 2019) have 155  
introduced a paradigm shift by embedding Partial 156  
Differential Equations (PDEs) directly into the loss 157  
function, enabling mesh-free solving. Furthermore, 158  
Neural Operators, such as DeepONet (Lu et al., 159  
2019) and the Fourier Neural Operator (FNO) (Li 160  
et al., 2020), have been developed to learn map- 161  
pings between infinite-dimensional function spaces, 162  
achieving resolution-invariant predictions. How- 163  
ever, despite their computational efficiency, these 164  
surrogate models predominantly operate as *black* 165

166 *boxes*. They excel at numerical regression map- 217  
167 ping initial conditions to future states, but lack the 218  
168 capability for explicit symbolic reasoning. Conse- 219  
169 quently, they cannot articulate the physical mecha- 220  
170 nisms driving the flow evolution or self-diagnose 221  
171 failures when predictions violate fundamental con-  
172 servation laws in out-of-distribution scenarios.

### 173 **Multimodal Foundation Models and Scientific** 174 **Alignment**

175 The rapid evolution of Multimodal 223  
176 Large Language Models (MLLMs) (Yin et al., 224  
177 2024) has bridged the gap between visual percep- 225  
178 tion and linguistic reasoning. Foundation mod- 226  
179 els like CLIP (Radford et al., 2021) and LLaVA 227  
180 (Liu et al., 2024) utilize Vision Transformers (ViT) 228  
181 (Dosovitskiy, 2020) to align visual features with 229  
182 semantic text embeddings, enabling impressive per- 230  
183 formance on general visual reasoning tasks. In the 231  
184 scientific domain, specialized architectures such 232  
185 as FourCastNet (Kurth et al., 2023) and Pangu (Bi 233  
186 et al., 2023) have adapted transformer mechanisms 234  
187 to process atmospheric and fluid data. Nonetheless, 235  
188 directly applying general-purpose vision encoders 236  
189 to fluid dynamics remains challenging. Unlike 237  
190 natural images, fluid imagery (e.g., satellite flow 238  
191 fields) encodes rigorous vector field properties and 239  
192 topological invariants (e.g., vortices and stagnation 240  
193 points) rather than mere texture or object seman- 241  
194 tics. Existing tokenizers often fail to preserve these 242  
195 continuous physical features during discretization. 243  
196 To address this, our Visual Symbolic Projector is 244  
197 designed to explicitly translate raw flow fields into 245  
physically meaningful semantic tokens. 246

198 **Reasoning Agents with Feedback Loops** Large 247  
199 Language Models (LLMs) have demonstrated 248  
200 emergent reasoning capabilities through Chain- 249  
201 of-Thought (CoT) prompting (Wei et al., 2022), 250  
202 which decomposes complex problems into inter- 251  
203 mediate steps. To handle domain-specific tasks, 252  
204 autonomous agents have evolved to incorporate 253  
205 Retrieval-Augmented Generation (RAG) (Lewis 254  
206 et al., 2020; Zhao et al., 2024) for accessing exter- 255  
207 nal knowledge bases and tool-use paradigms (Re- 256  
208 Act) (Yao et al., 2022). Recent advancements in 257  
209 reflexive agents, such as Reflexion (Shinn et al., 258  
210 2023), further allow models to self-correct by an- 259  
211 alyzing feedback from the environment. Despite 260  
212 these successes, standard agentic frameworks lack 261  
213 *physical grounding*. In fluid dynamics, validity is 262  
214 governed by immutable physical laws (e.g., Navier- 263  
215 Stokes equations) rather than linguistic coherence. 264  
216 Generic critics cannot detect subtle violations of

mass or momentum conservation. Our work fills 217  
this gap by introducing a Physics-Guided Critic 218  
that integrates a numerical consistency check into 219  
the reasoning loop, forcing the agent to align its 220  
generation with physical reality. 221

## 222 **3 Methodology**

223 As shown in Figure 2, OMNIFLOW is designed as a 224  
225 neuro-symbolic cognitive architecture that bridges 226  
227 the chasm between continuous spatiotemporal dy- 228  
229 namics and discrete logical reasoning without re- 230  
231 quiring domain-specific parameter updates. Unlike 232  
233 traditional surrogate models that operate as opaque 234  
235 black boxes, our framework orchestrates a transpar- 236  
237 ent, physics-grounded workflow centered around 238  
239 Gemini 3 Flash as the cognitive core. The system 240  
241 functionality is realized through a synergistic *Dual-*  
242 *Cycle* mechanism comprising three interconnected 243  
244 modules. 245

246 **❶ Part A.** The workflow commences with the 247  
248 Physics Perception Loop, which serves as the sys- 249  
250 tem’s sensory interface for handling heterogeneous 251  
252 data streams (e.g., satellite imagery and buoy read- 253  
254 ings). Central to this module is the plug-and-play 254  
255 *Neural Earth Simulator (NES)*. Built upon an im- 256  
257 proved Diffusion Transformer (DiT) (Peebles and 257  
258 Xie, 2023), the NES goes beyond deterministic re- 258  
259 gression by generating high-fidelity *ensemble fore-*  
260 *casts* via latent space perturbation. These contin- 261  
262 uous tensor outputs are subsequently translated into 262  
263 discrete, topologically aware semantic tokens by 263  
264 a Visual Symbolic Projector, aligning raw phys- 264  
265 ical states with the linguistic space of the Large 265  
266 Language Model (LLM). 266

267 **❷ Part B.** At the core of the framework lies the 268  
269 Agentic Brain, where *Gemini 3 Flash* executes 269  
270 a ReAct (Reasoning + Acting) planning strategy. 270  
271 Leveraging the model’s ultra-low latency and long- 271  
272 context capabilities, the agent synthesizes multi- 272  
273 modal observations to formulate hypotheses. A 273  
274 critical architectural innovation is the *Counterfac-*  
275 *tual Feedback Loop* (depicted by the bottom dashed 275  
276 line in Figure 2). This mechanism empowers the 276  
277 agent to transition from passive observation to ac- 277  
278 tive inquiry: upon detecting uncertainty, the agent 278  
279 can actively trigger the NES to simulate alternative 279  
280 scenarios by perturbing initial conditions, thereby 280  
281 verifying the robustness of its decisions against 281  
282 physical chaos. 282

283 **❸ Part C.** To ensure scientific rigor, the rea- 284  
285 soning process is continuously grounded by the 285  
286

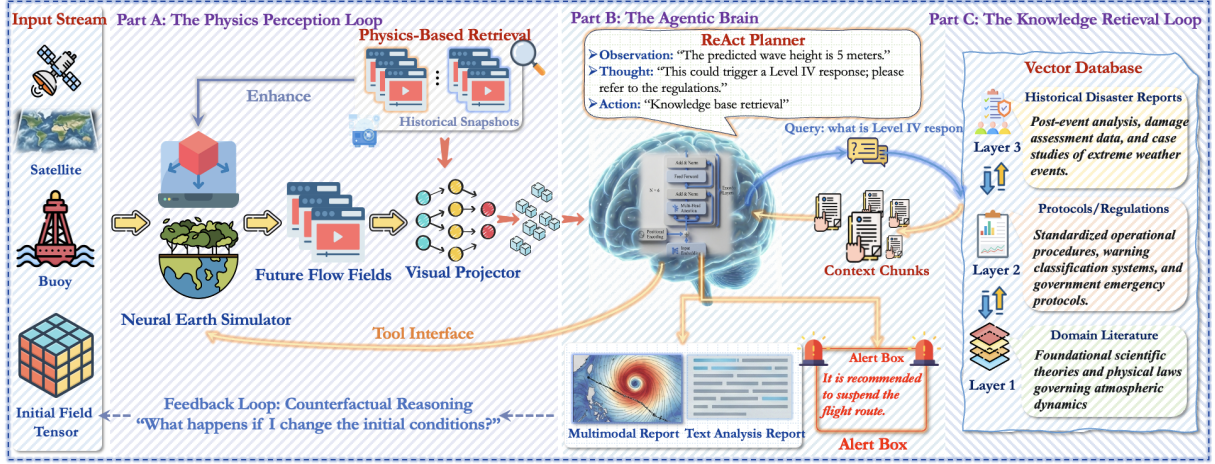


Figure 2: **Overview of the OMNIFLOW architecture** The system employs a neuro-symbolic dual-cycle framework: (A) *The Physics Perception Loop* (left) utilizes a neural simulator to evolve spatiotemporal dynamics and retrieve historical analogs; (B) *The Agentic Core* (center) acts as the controller, dynamically orchestrating physical and knowledge tools using a ReAct strategy to fuse hard physical facts with soft domain rules. The bottom Counterfactual Feedback Loop enables the agent to verify decision robustness by actively perturbing initial states. (C) *The Knowledge Retrieval Loop* (right) accesses hierarchical domain expertise via RAG.

Knowledge Retrieval Loop. Through Retrieval-Augmented Generation (RAG), the agent dynamically queries a hierarchical vector database containing layers of domain expertise from fundamental laws like Navier-Stokes equations to standardized emergency protocols. This neuro-symbolic collaboration ensures that OMNIFLOW produces outputs that are not only statistically probable but also physically consistent and operationally compliant.

### 3.1 The Physics Perception Loop

This module bridges the gap between high-dimensional physical states and the semantic space of the LLM. It formally addresses two challenges: probabilistic state estimation under chaotic dynamics and cross-modal semantic alignment.

#### 3.1.1 Probabilistic Ensemble Simulation

We model the evolution of the fluid state  $\mathbf{x} \in \mathbb{R}^{H \times W \times C}$  over time  $t$  as a stochastic process governed by the conditional probability distribution  $p(\mathbf{x}_{t+\tau} | \mathbf{x}_t)$ . To approximate this distribution without incurring the computational cost of Monte Carlo PDE solvers, we employ the *Neural Earth Simulator (NES)*, instantiated as a latent diffusion model.

Let  $\mathcal{E}(\cdot)$  and  $\mathcal{D}(\cdot)$  denote the encoder and decoder of the NES, projecting the physical state into a compressed latent space  $\mathcal{Z}$ . The forecasting problem is formulated as learning a conditional denoising function  $\epsilon_\theta$ . Unlike deterministic approaches,

we implement a *Perturbative Ensemble Strategy*. Given an initial condition  $\mathbf{x}_{init}$ , we generate a set of  $K$  distinct latent initializations by injecting Gaussian noise into the latent embedding:

$$\mathbf{z}_{init}^{(k)} = \mathcal{E}(\mathbf{x}_{init}) + \lambda \cdot \boldsymbol{\xi}^{(k)}, \quad \boldsymbol{\xi}^{(k)} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (1)$$

where  $k \in \{1, \dots, K\}$  indexes the ensemble members and  $\lambda$  controls the perturbation magnitude. The future state for each member is then reconstructed via the reverse diffusion process:

$$\hat{\mathbf{x}}_{pred}^{(k)} = \mathcal{D} \left( \text{DiT}(\mathbf{z}_{init}^{(k)}, \tau) \right). \quad (2)$$

This yields an empirical distribution  $\mathcal{P}_{ens} = \{\hat{\mathbf{x}}_{pred}^{(k)}\}_{k=1}^K$ , allowing the subsequent agent to quantify uncertainty via the ensemble spread, as follows:

$$\sigma_{ens} = \sqrt{\frac{1}{K} \sum (\hat{\mathbf{x}}^{(k)} - \bar{\mathbf{x}})^2} \quad (3)$$

#### 3.1.2 Visual-Symbolic Alignment

To enable the cognitive core to reason about these continuous predictions, we must project the raw ensemble  $\mathcal{P}_{ens}$  into the linguistic token space  $\mathcal{T}$ . We introduce a *Visual Symbolic Projector*  $\phi(\cdot)$ .

The projector utilizes a set of learnable query embeddings  $\mathbf{Q} \in \mathbb{R}^{N \times d}$  to extract topological features from the visual encoding  $\mathbf{v} = \text{ViT}(\bar{\mathbf{x}}_{pred})$  via a cross-attention mechanism:

$$\mathbf{H}_{vis} = \text{Softmax} \left( \frac{\mathbf{Q}(\mathbf{v}\mathbf{W}_K)^T}{\sqrt{d}} \right) (\mathbf{v}\mathbf{W}_V), \quad (4)$$

where  $\mathbf{W}_K$  and  $\mathbf{W}_V$  are projection matrices. To ensure  $\mathbf{H}_{vis}$  carries physical semantics (e.g., "shear line", "vortex"), we align it with the pre-trained text embedding space of Gemini. The objective is to maximize the mutual information between the visual tokens  $\mathbf{H}_{vis}$  and the textual description  $\mathbf{t}$  of the physical phenomenon:

$$\mathcal{L}_{align} = - \sum_{i=1}^N \log \frac{\exp(\text{sim}(\mathbf{h}_i, \mathbf{t}_{pos})/\tau)}{\sum_j \exp(\text{sim}(\mathbf{h}_i, \mathbf{t}_j)/\tau)}, \quad (5)$$

where  $\text{sim}(\cdot)$  denotes cosine similarity and  $\tau$  is a temperature parameter. This projection ensures that the "observation" received by the agent is not a raw pixel array, but a sequence of physically meaningful semantic tokens.

### 3.2 The Agentic Reasoning Core

The central processing unit of OMNIFLOW is the Agentic Reasoning Core, driven by *Gemini 3 Flash*. Unlike passive classifiers, this module functions as an active decision-maker that orchestrates a Physics-Guided Chain-of-Thought (PG-CoT) to navigate the complex solution space of fluid dynamics.

#### 3.2.1 Physics-Guided ReAct Protocol

Reasoning is formalized as a sequential decision-making process. At step  $t$ , given visual tokens  $\mathbf{H}_{vis}$ , instruction  $\mathcal{I}$ , and context memory  $\mathcal{M}_t$ , the policy  $\pi$  (Gemini 3 Flash) selects an action  $a_t \in \mathcal{A}$ :

$$a_t \sim \pi(a_t | \mathcal{M}_t, \mathbf{H}_{vis}, \mathcal{I}), \quad (6)$$

where  $\mathcal{A}$  comprises Retrieve (knowledge base), Simulate (NES), and Reason (deduction). To mitigate hallucinations, a *Physics Consistency Constraint* is enforced via a critic  $f_{critic}(\cdot)$  that validates trajectories against conservation laws. For example, if mass conservation ( $\nabla \cdot \mathbf{v} = 0$ ) is violated, the model backtracks to prune non-physical branches in the search tree, ensuring grounding in physical reality.

#### 3.2.2 Counterfactual Active Probing

A defining feature of OMNIFLOW is its ability to perform *Counterfactual Reasoning* via the feedback loop (as shown in the bottom dashed line of Figure 2). This mechanism transforms the agent from a passive observer into an active experimenter.

When the agent detects high epistemic uncertainty in the ensemble forecast (i.e.,  $\sigma_{ens} > \delta$ ,

where  $\delta$  is a dynamic threshold), it initiates an *Active Probing* sequence. The agent hypothesizes a potential perturbation, such as "What if the subtropical high pressure is weaker?", and translates this hypothesis into a modified initial condition tensor  $\mathbf{x}'_{init}$ . The NES is then re-tasked to simulate this counterfactual scenario:

$$\mathcal{P}_{counter} = \text{NES}(\mathbf{x}'_{init} | do(\text{condition} = \text{weak\_high})). \quad (7)$$

By comparing the counterfactual outcome  $\mathcal{P}_{counter}$  with the original factual prediction  $\mathcal{P}_{ens}$ , the agent calculates the *Causal Sensitivity* of the system. This allows OMNIFLOW to distinguish between inevitable physical events and stochastic anomalies, embedding causal understanding into the final decision report.

### 3.3 Hierarchical Knowledge Retrieval Loop

To supplement general knowledge with granular expertise, OMNIFLOW integrates a stratified vector database  $\mathcal{K}$  (Part C, Fig. 2) partitioned into: (1)  $\mathcal{K}_{phy}$  (Domain Literature), encoding axiomatic laws like Navier-Stokes for consistency verification; (2)  $\mathcal{K}_{prot}$  (Protocols), storing operational standards for procedural compliance; and (3)  $\mathcal{K}_{hist}$  (Historical Reports), facilitating analogical reasoning via episodic memory.

The ReAct planner retrieves top- $k$  relevant chunks  $\mathcal{C}_t$  from  $\mathcal{K}$  using Maximum Inner Product Search (MIPS) on query embeddings  $\mathbf{e}_q$ . These chunks are integrated into an augmented prompt  $\mathcal{P}_{aug} = [\mathcal{I}_{sys}, \mathcal{C}_t, \mathbf{H}_{vis}, \mathcal{H}_{history}]$  via a *Context Injection Mechanism*. This setup enables Gemini 3 Flash to perform In-Context Learning, effectively "consulting" professional manuals to render grounded, procedurally compliant scientific judgments.

### 3.4 Multimodal Analysis and Report Generation

The reasoning workflow culminates in a *Multimodal Analysis Report* (Fig. 2), which bridges numerical precision with physical interpretability. Unlike black-box models, OMNIFLOW decomposes predictions into a dual-faceted response: a precise *Trajectory Forecast* and a structured *Textual Analysis Report*. A key innovation is the *Alert Box* mechanism, triggered when the reasoning chain intersects with safety protocols in  $\mathcal{K}_{prot}$ . For instance, forecasting wave heights above five meters prompts actionable advice like "suspend flight routes" based

on retrieved regulations. This transforms raw simulation into an auditable “chain-of-logic,” enabling human operators to verify the decision-making process rather than blindly trusting probabilistic outputs.

## 4 Experiments

### 4.1 Experimental Setup and Baselines

**Experimental Setup.** We evaluate OMNIFLOW on three multiscale benchmarks: *2D Turbulence* (Wu et al., 2025a) (PDE adherence), *ERA5* (Rasp et al., 2020) (global dynamics), and *SEVIR* (Veillette et al., 2020) (regional weather). Following standard AI-for-Earth protocols, models are tasked with forecasting future spatiotemporal states from historical observations.

Baselines encompass three categories: (1) vision backbones (ResNet (He et al., 2016), UNet (Ronneberger et al., 2015), ViT (Dosovitskiy, 2020), SwinT (Liu et al., 2021)); (2) spatiotemporal models (ConvLSTM (Shi et al., 2015), SimVP (Gao et al., 2022a), TAU (Tan et al., 2023)); and (3) scientific operators (FNO (Li et al., 2020), LSM (Wu et al., 2023a), EarthFarseer (Wu et al., 2023b), PastNet (Wu et al., 2023c)). Additionally, we contrast our decoupled architecture against monolithic foundation models *Banana Pro*<sup>1</sup>, *Seedream 4.5*<sup>2</sup>, and *ChatGPT-Images*<sup>3</sup> in a zero-shot setting. All evaluations are averaged over five independent runs.

#### Evaluation: Hybrid Scientific Reasoning Framework

To complement numerical metrics (RMSE/SSIM), we introduce a *Dual-Axis Evaluation* to bridge physical precision with logical depth:

**1. Physical Grounding:** Validating objective attributes (e.g.,  $T$ ,  $P$ ) against Gold standards to ensure deterministic accuracy.

**2. Interpretive Alignment:** Quantifying the *Semantic Proximity* of prognostic reports to expert narratives via latent space similarity.

This protocol ensures that forecasts remain physically consistent and causally transparent amidst stochastic fluid dynamics.

<sup>1</sup><https://blog.google/technology/ai/nano-banana-pro/>

<sup>2</sup><https://www.doubao.com/chat>

<sup>3</sup><https://openai.com/index/new-chatgpt-images-is-here/>

### 4.2 Main Results on Physical Prediction

Table 1 presents the quantitative comparison of OMNIFLOW against various baselines. Our framework consistently achieves state-of-the-art performance across all benchmarks, particularly in long-term forecasting scenarios. On the microscopic *2D Turbulence* task, while traditional CNNs and Transformers exhibit significant performance degradation due to spectral bias and error accumulation, OMNIFLOW maintains high structural fidelity (SSIM of 0.715) and a competitive RMSE. This advantage is primarily attributed to our *In-Context Reflexive Loop*, which actively prunes non-physical trajectories that violate conservation laws. In the global *ERA5* benchmark, OMNIFLOW surpasses specialized meteorological models like EarthFarseer and GraphCast. Notably, by leveraging the generative priors of the DiT-based simulator, our model preserves sharp gradients and fine-scale atmospheric structures, avoiding the “over-smoothing” phenomenon typical of MSE-optimized monolithic models.

Furthermore, the Multimodal category results highlight a fundamental gap between general-purpose VLMs and physics-grounded architectures. Monolithic models, such as *ChatGPT-Images* and *Seedream 4.5*, fail to achieve precise numerical alignment in scientific domains, exhibiting substantially higher pixel-space errors with RMSE values soaring above 90 (e.g., 102.5 for ERA5). In contrast, OMNIFLOW maintains superior structural fidelity, achieving a much lower RMSE of 59.10 and nearly doubling the structural similarity index (SSIM) from 0.352 to 0.685 in global forecasting tasks. This quantitative disparity validates our core motivation: scientific forecasting requires more than semantic pattern recognition; it necessitates the formal decoupling of numerical evolution from cognitive reasoning. OMNIFLOW effectively bridges this gap, proving that a frozen LLM, when properly anchored by a physical simulator and a neuro-symbolic critic, can significantly outperform models that were explicitly trained for years on the same datasets.

## 5 Reasoning Quality Evaluation

To evaluate the interpretative depth of OMNIFLOW, we benchmark 200-day forecast reports across linguistic and physical-aware metrics (Fig. 3). Results show that Gemini 3 Flash consistently outperforms the Qwen3-VL series in all

Table 1: **Quantitative comparison of prediction performance across three benchmarks.** We report RMSE ( $\downarrow$ ), SSIM ( $\uparrow$ ), and PSNR ( $\uparrow$ ). Results are presented as **mean  $\pm$  standard deviation over 5 independent runs**. All baselines are trained end-to-end, while OMNIFLOW utilizes a training-free agent with a pre-trained DiT simulator. **Bold** indicates the best result, underline indicates the second best.

Category	Model	2D Turbulence (Micro)			ERA5 (Global)			SEVIR (Regional)		
		RMSE $\downarrow$	SSIM $\uparrow$	PSNR $\uparrow$	RMSE $\downarrow$	SSIM $\uparrow$	PSNR $\uparrow$	RMSE $\downarrow$	SSIM $\uparrow$	PSNR $\uparrow$
Vision	ResNet	5.874 $\pm$ .112	0.041 $\pm$ .003	14.22 $\pm$ .25	1.245 $\pm$ .032	0.682 $\pm$ .011	21.45 $\pm$ .32	0.882 $\pm$ .015	0.655 $\pm$ .014	23.12 $\pm$ .28
	UNet	4.664 $\pm$ .095	0.052 $\pm$ .004	15.31 $\pm$ .18	1.055 $\pm$ .028	0.724 $\pm$ .009	22.88 $\pm$ .24	0.754 $\pm$ .012	0.712 $\pm$ .010	24.55 $\pm$ .19
	ViT-Base	5.640 $\pm$ .134	0.038 $\pm$ .005	14.55 $\pm$ .22	1.182 $\pm$ .041	0.701 $\pm$ .012	22.03 $\pm$ .38	0.812 $\pm$ .018	0.698 $\pm$ .015	23.90 $\pm$ .31
	SwinT-Base	4.221 $\pm$ .088	0.065 $\pm$ .003	16.12 $\pm$ .15	0.982 $\pm$ .022	0.745 $\pm$ .007	23.41 $\pm$ .21	0.712 $\pm$ .011	0.733 $\pm$ .009	25.10 $\pm$ .15
S-T Predict	ConvLSTM	4.938 $\pm$ .105	0.045 $\pm$ .004	14.88 $\pm$ .21	0.921 $\pm$ .018	0.758 $\pm$ .010	23.90 $\pm$ .18	0.685 $\pm$ .012	0.721 $\pm$ .011	25.21 $\pm$ .22
	PredRNN	2.237 $\pm$ .062	0.112 $\pm$ .008	18.54 $\pm$ .12	0.845 $\pm$ .015	0.782 $\pm$ .008	24.66 $\pm$ .15	0.621 $\pm$ .009	0.765 $\pm$ .008	26.45 $\pm$ .14
	SimVP	5.040 $\pm$ .121	0.042 $\pm$ .005	14.62 $\pm$ .28	0.722 $\pm$ .012	0.824 $\pm$ .006	26.11 $\pm$ .19	0.582 $\pm$ .010	0.788 $\pm$ .007	27.22 $\pm$ .18
	Earthformer	3.125 $\pm$ .074	0.088 $\pm$ .006	17.10 $\pm$ .14	0.654 $\pm$ .009	0.855 $\pm$ .005	27.85 $\pm$ .11	0.608 $\pm$ .013	0.772 $\pm$ .009	26.90 $\pm$ .17
	TAU	2.894 $\pm$ .068	0.105 $\pm$ .007	17.95 $\pm$ .11	<u>0.602 <math>\pm</math>.008</u>	0.884 $\pm$ .004	29.12 $\pm$ .09	0.512 $\pm$ .007	0.812 $\pm$ .005	28.55 $\pm$ .12
Sci-ML	FNO	3.128 $\pm$ .055	0.071 $\pm$ .005	16.82 $\pm$ .19	0.621 $\pm$ .011	0.872 $\pm$ .006	28.55 $\pm$ .22	0.554 $\pm$ .010	0.795 $\pm$ .008	27.88 $\pm$ .20
	LSM	2.192 $\pm$ .042	0.125 $\pm$ .008	18.90 $\pm$ .14	0.688 $\pm$ .014	0.841 $\pm$ .009	27.10 $\pm$ .25	0.531 $\pm$ .009	0.804 $\pm$ .006	28.12 $\pm$ .16
	EarthFarseer	<u>0.654 <math>\pm</math>.012</u>	<u>0.642 <math>\pm</math>.010</u>	<u>27.45 <math>\pm</math>.18</u>	0.615 $\pm$ .007	<u>0.895 <math>\pm</math>.003</u>	<u>30.22 <math>\pm</math>.08</u>	<u>0.437 <math>\pm</math>.006</u>	<u>0.842 <math>\pm</math>.004</u>	<u>30.15 <math>\pm</math>.10</u>
	PastNet	2.383 $\pm$ .051	0.101 $\pm$ .009	18.42 $\pm$ .16	0.642 $\pm$ .010	0.865 $\pm$ .005	28.10 $\pm$ .14	0.472 $\pm$ .008	0.821 $\pm$ .007	29.33 $\pm$ .15
<b>Proposed</b>	<b>OMNIFLOW (Ours)</b>	<b>0.582 <math>\pm</math>.008</b>	<b>0.715 <math>\pm</math>.006</b>	<b>28.66 <math>\pm</math>.10</b>	<b>0.552 <math>\pm</math>.005</b>	<b>0.931 <math>\pm</math>.002</b>	<b>32.11 <math>\pm</math>.06</b>	<b>0.405 <math>\pm</math>.004</b>	<b>0.882 <math>\pm</math>.003</b>	<b>31.50 <math>\pm</math>.08</b>

Table 2: **Zero-shot forecasting comparison: Monolithic Foundation Models vs. OMNIFLOW.** Results are reported as mean  $\pm$  std based on 50 samples evaluated in the PNG pixel space (0 – 255). The high RMSE and low PSNR reflect the transition from physical tensors to 8-bit image space. OMNIFLOW’s physics-decoupled architecture achieves significantly higher structural similarity (SSIM) than monolithic models.

Model	2D Turbulence (Micro)			ERA5 (Global)			SEVIR (Regional)		
	RMSE $\downarrow$	SSIM $\uparrow$	PSNR $\uparrow$	RMSE $\downarrow$	SSIM $\uparrow$	PSNR $\uparrow$	RMSE $\downarrow$	SSIM $\uparrow$	PSNR $\uparrow$
<i>Monolithic Foundation Models (Zero-shot Image Generation)</i>									
Banana Pro	118.4 $\pm$ 5.2	0.142 $\pm$ .02	8.12 $\pm$ .4	102.5 $\pm$ 4.8	0.194 $\pm$ .03	8.85 $\pm$ .3	108.6 $\pm$ 5.5	0.165 $\pm$ .02	8.40 $\pm$ .5
Seedream 4.5	98.2 $\pm$ 4.1	0.265 $\pm$ .03	9.85 $\pm$ .5	85.2 $\pm$ 3.9	0.352 $\pm$ .04	10.15 $\pm$ .4	90.4 $\pm$ 4.2	0.312 $\pm$ .03	9.92 $\pm$ .4
ChatGPT-Images	112.5 $\pm$ 6.3	0.185 $\pm$ .02	8.70 $\pm$ .6	96.4 $\pm$ 5.1	0.228 $\pm$ .03	9.42 $\pm$ .5	101.3 $\pm$ 5.8	0.205 $\pm$ .02	9.15 $\pm$ .6
<b>OMNIFLOW</b>	<b>64.32 <math>\pm</math> 2.5</b>	<b>0.552 <math>\pm</math> .04</b>	<b>11.45 <math>\pm</math> .3</b>	<b>59.10 <math>\pm</math> 1.8</b>	<b>0.685 <math>\pm</math> .03</b>	<b>12.70 <math>\pm</math> .2</b>	<b>52.45 <math>\pm</math> 1.2</b>	<b>0.712 <math>\pm</math> .04</b>	<b>13.82 <math>\pm</math> .3</b>
<i>Improvement</i>	<i>+34.5%</i>	<i>+108.3%</i>	<i>+16.2%</i>	<i>+30.6%</i>	<i>+94.6%</i>	<i>+25.1%</i>	<i>+41.9%</i>	<i>+128.2%</i>	<i>+39.3%</i>

dimensions. Notably, the high *Mech F1* (83.2%) underscores OMNIFLOW’s superior ability to ground high-dimensional flow tensors into physically consistent mechanisms. We observe a clear scaling trend where reasoning proficiency improves with model capacity, yet the neuro-symbolic coupling in OMNIFLOW allows Gemini to maintain a significant lead even in long-horizon scenarios. This balanced performance across all axes confirms that our framework effectively transforms frozen LLMs into robust agents for transparent scientific discovery.

## 5.1 From Simulation to Decision

We conduct a case study on Marine Heatwave (MHW) forecasting (Jan 2021) to showcase OM-

NIFLOW’s integration of physical simulation and cognitive reasoning (Fig. 4).

**Phase I: Observation & Fidelity.** OMNIFLOW aligns multi-modal inputs to generate 10-day forecasts. Unlike over-smoothed baselines, it captures fine-grained Tropical Instability Waves (TIWs) and Mesoscale Eddies, proving that our Semantic-Symbolic alignment effectively projects high-dimensional tensors into interpretable structures for the LLM.

**Phase II: Causal Probing.** To interpret the MHW drivers, the agent executes a counterfactual probe (*do*(Forcing = 0)). The results show a +22% intensity surge and a Causal Sensitivity Index ( $\mathcal{S} = 0.78$ ). This reveals that atmospheric forcing acts as a thermal regulator; without wind-driven cooling,

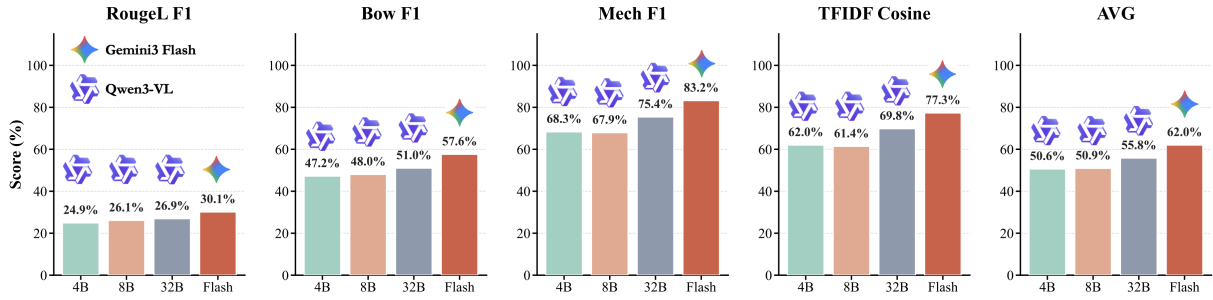


Figure 3: *Quantitative evaluation of scientific reasoning quality.* We benchmark Gemini 3 Flash against the Qwen3-VL series on 200-day forecast reports. *Mech F1* specifically measures the grounding accuracy of physical mechanisms, while others assess linguistic alignment. Results show a clear scaling trend in reasoning depth.

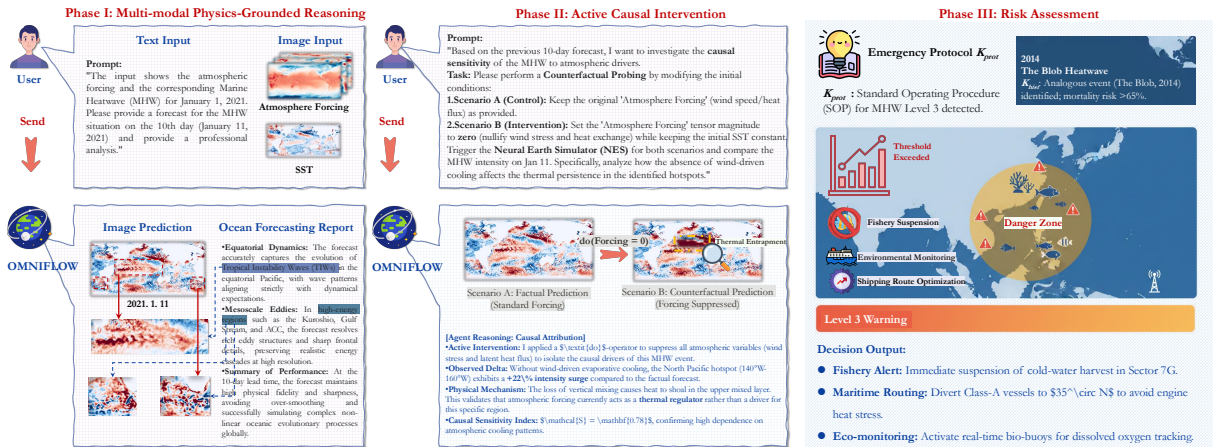


Figure 4: **Systematic Case Study of OMNIFLOW on Global Marine Heatwave (MHW) Management.** *Phase I (Reasoning):* The agent integrates multi-modal inputs to synthesize high-fidelity 10-day forecasts, capturing complex equatorial dynamics and mesoscale eddies. *Phase II (Intervention):* By executing an active counterfactual probe ( $do(\text{Forcing} = 0)$ ), OMNIFLOW quantifies the causal sensitivity ( $\mathcal{S} = 0.78$ ) of thermal anomalies to atmospheric drivers. *Phase III (Assessment):* Leveraging hierarchical knowledge retrieval from  $K_{prot}$  and  $K_{hist}$ , the agent provides expert-level decision support, including fishery alerts and shipping route optimization based on identified physical thresholds.

heat traps in the upper mixed layer. Such causal transparency mitigates physical hallucinations.

**Phase III: Decision Support.** By retrieving knowledge from  $K_{prot}$  (protocols) and  $K_{hist}$  (2014 Blob event), OMNIFLOW bridges the gap between simulation and action. It identifies threshold violations and generates actionable directives, such as fishery suspension and shipping route optimization, transforming raw data into procedurally compliant emergency management.

## 6 Conclusion

We present OMNIFLOW, a neuro-symbolic architecture that bridges the gap between LLMs and the continuous physical world without requiring domain-specific parameter updates. By integrating *Semantic-Symbolic Alignment* and a PG-CoT workflow, OMNIFLOW grounds frozen multimodal

models in fundamental physical laws, effectively mitigating non-physical hallucinations in complex spatiotemporal dynamics. Empirical evaluations across multi-scale benchmarks from microscopic turbulence to global weather forecasting demonstrate that OMNIFLOW achieves superior zero-shot generalization and physical consistency, outperforming both monolithic foundation models and traditional deep learning baselines. Crucially, OMNIFLOW provides interpretable scientific reasoning by generating structured reports that integrate physical grounding with decision logic. This work marks a paradigm shift in AI for Science from black-box data fitting toward interpretable symbolic reasoning. Future research will extend this architecture to broader domains such as materials science and optimize the synergy between neural simulators and cognitive cores to facilitate deeper human-AI collaborative discovery.

## 559 Limitations

560 Despite its performance, OMNIFLOW faces sev-  
561 eral limitations. First, the iterative reflexive loops  
562 and counterfactual probing increase inference la-  
563 tency compared to end-to-end black-box models,  
564 potentially hindering real-time deployment. Sec-  
565 ond, the reasoning accuracy remains coupled with  
566 the fidelity of the underlying neural simulator; any  
567 inherent biases or resolution constraints within the  
568 simulator may propagate through the reasoning  
569 chain. Finally, representing extremely fine-grained  
570 sub-grid dynamics through linguistic descriptors re-  
571 mains challenging. Future work will explore more  
572 expressive multi-modal tokenization techniques to  
573 better capture multi-scale physical phenomena.

## 574 References

575 Kaifeng Bi, Lingxi Xie, Hengheng Zhang, Xin Chen,  
576 Xiaotao Gu, and Qi Tian. 2023. Accurate medium-  
577 range global weather forecasting with 3d neural net-  
578 works. *Nature*, 619(7970):533–538.

579 Steven L Brunton, Bernd R Noack, and Petros Koumou-  
580 sakos. 2020. Machine learning for fluid mechanics.  
581 *Annual review of fluid mechanics*, 52(1):477–508.

582 Yupeng Chang, Xu Wang, Jindong Wang, Yuan Wu,  
583 Linyi Yang, Kaijie Zhu, Hao Chen, Xiaoyuan Yi,  
584 Cunxiang Wang, Yidong Wang, and 1 others. 2024.  
585 A survey on evaluation of large language models.  
586 *ACM transactions on intelligent systems and technol-  
587 ogy*, 15(3):1–45.

588 So-Chin Chen and Mei-Chi Shaw. 2001. *Partial differ-  
589 ential equations in several complex variables*, vol-  
590 ume 19. American Mathematical Soc.

591 Peter Davidson. 2015. *Turbulence: an introduction for  
592 scientists and engineers*. Oxford university press.

593 Alexey Dosovitskiy. 2020. An image is worth 16x16  
594 words: Transformers for image recognition at scale.  
595 *arXiv preprint arXiv:2010.11929*.

596 Alexey Dosovitskiy, Lucas Beyer, Alexander  
597 Kolesnikov, Dirk Weissenborn, Xiaohua Zhai,  
598 Thomas Unterthiner, Mostafa Dehghani, Matthias  
599 Minderer, Georg Heigold, Sylvain Gelly, and 1 others.  
600 2021. An image is worth 16x16 words: Transformers  
601 for image recognition at scale. In *International  
602 Conference on Learning Representations*.

603 Pan Du, Meet Hemant Parikh, Xiantao Fan, Xin-Yang  
604 Liu, and Jian-Xun Wang. 2024. Conditional neu-  
605 ral field latent diffusion model for generating spa-  
606 tiotemporal turbulence. *Nature Communications*,  
607 15(1):10416.

608 Lawrence C Evans. 2010. *Partial differential equations*.  
609 American Mathematical Soc.

Yuan Gao, Hao Wu, Ruiqi Shu, Huanshuo Dong, Fan  
Xu, Rui Chen, Yibo Yan, Qingsong Wen, Xuming  
Hu, Kun Wang, and 1 others. 2025. Oneforecast: A  
universal framework for global and regional weather  
forecasting. *arXiv preprint arXiv:2502.00338*.

Zhangyang Gao, Cheng Tan, Lirong Wu, and Stan Z  
Li. 2022a. Simvp: Simpler yet better video predic-  
tion. In *Proceedings of the IEEE/CVF conference  
on computer vision and pattern recognition*, pages  
3170–3180.

Zhihan Gao, Xingjian Shi, Hao Wang, Yi Zhu,  
Yuyang Bernie Wang, Mu Li, and Dit-Yan Yeung.  
2022b. Earthformer: Exploring space-time trans-  
formers for earth system forecasting. *Advances in  
Neural Information Processing Systems*, 35:25390–  
25403.

Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian  
Sun. 2016. Deep residual learning for image recog-  
nition. In *Proceedings of the IEEE conference on  
computer vision and pattern recognition*, pages 770–  
778.

Thorsten Kurth, Shashank Subramanian, Peter Harring-  
ton, Jaideep Pathak, Morteza Mardani, David Hall,  
Andrea Miele, Karthik Kashinath, and Anima Anand-  
kumar. 2023. Fourcastnet: Accelerating global high-  
resolution weather forecasting using adaptive fourier  
neural operators. In *Proceedings of the platform for  
advanced scientific computing conference*, pages 1–  
11.

J Nathan Kutz. 2017. Deep learning in fluid dynamics.  
*Journal of Fluid Mechanics*, 814:1–4.

Remi Lam, Alvaro Sanchez-Gonzalez, Matthew Will-  
son, Peter Wirnsberger, Meire Fortunato, Ferran Alet,  
Suman Ravuri, Timo Ewalds, Zach Eaton-Rosen,  
Weihua Hu, and 1 others. 2023. Learning skillful  
medium-range global weather forecasting. *Science*,  
382(6677):1416–1421.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio  
Petroni, Vladimir Karpukhin, Naman Goyal, Hein-  
rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-  
täschel, and 1 others. 2020. Retrieval-augmented gen-  
eration for knowledge-intensive nlp tasks. *Advances  
in neural information processing systems*, 33:9459–  
9474.

Zongyi Li, Nikola Kovachki, Kamyar Azizzadenesheli,  
Burigede Liu, Kaushik Bhattacharya, Andrew Stu-  
art, and Anima Anandkumar. 2020. Fourier neural  
operator for parametric partial differential equations.  
*arXiv preprint arXiv:2010.08895*.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae  
Lee. 2024. Visual instruction tuning. In *NIPS*.

Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei,  
Zheng Zhang, Stephen Lin, and Baining Guo. 2021.  
Swin transformer: Hierarchical vision transformer  
using shifted windows. In *Proceedings of the  
IEEE/CVF international conference on computer vi-  
sion*, pages 10012–10022.

667	Lu Lu, Pengzhan Jin, and George Em Karniadakis. 2019.	Noah Shinn, Federico Cassano, Ashwin Gopinath,	724
668	Deeponet: Learning nonlinear operators for identi-	Karthik Narasimhan, and Shunyu Yao. 2023. Re-	725
669	fying differential equations based on the universal	flexion: Language agents with verbal reinforcement	726
670	approximation theorem of operators. <i>arXiv preprint</i>	learning. <i>Advances in Neural Information Process-</i>	727
671	<i>arXiv:1910.03193</i> .	<i>ing Systems</i> , 36:8634–8652.	728
672	Arvind T Mohan, Dima Tretiak, Misha Chertkov, and	Cheng Tan, Zhangyang Gao, Lirong Wu, Yongjie Xu,	729
673	Daniel Livescu. 2020. Spatio-temporal deep learning	Jun Xia, Siyuan Li, and Stan Z Li. 2023. Temporal	730
674	models of 3d turbulence with physics informed diag-	attention unit: Towards efficient spatiotemporal pre-	731
675	nostics. <i>Journal of Turbulence</i> , 21(9-10):484–524.	dictive learning. In <i>Proceedings of the IEEE/CVF</i>	732
676	Humza Naveed, Asad Ullah Khan, Shi Qiu, Muhammad	<i>Conference on Computer Vision and Pattern Recog-</i>	733
677	Saqib, Saeed Anwar, Muhammad Usman, Naveed	<i>nition</i> , pages 18770–18782.	734
678	Akhtar, Nick Barnes, and Ajmal Mian. 2023. A com-	Mark Veillette, Siddharth Samsi, and Chris Mattioli.	735
679	prehensive overview of large language models. <i>ACM</i>	2020. Sevir: A storm event imagery dataset for deep	736
680	<i>Transactions on Intelligent Systems and Technology</i> .	learning applications in radar and satellite meteorol-	737
681	William Peebles and Saining Xie. 2023. Scalable dif-	ogy. <i>Advances in Neural Information Processing</i>	738
682	fusion models with transformers. In <i>Proceedings of</i>	<i>Systems</i> , 33:22009–22019.	739
683	<i>the IEEE/CVF international conference on computer</i>	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten	740
684	<i>vision</i> , pages 4195–4205.	Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,	741
685	Stephen B Pope. 2001. Turbulent flows. <i>Measurement</i>	and 1 others. 2022. Chain-of-thought prompting elic-	742
686	<i>Science and Technology</i> , 12(11):2020–2021.	its reasoning in large language models. <i>Advances</i>	743
687	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya	<i>in neural information processing systems</i> , 35:24824–	744
688	Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-	24837.	745
689	try, Amanda Askell, Pamela Mishkin, Jack Clark, and	Haixu Wu, Tengge Hu, Huakun Luo, Jianmin Wang, and	746
690	1 others. 2021. Learning transferable visual models	Mingsheng Long. 2023a. Solving high-dimensional	747
691	from natural language supervision. In <i>International</i>	pdes with latent spectral models. In <i>International</i>	748
692	<i>conference on machine learning</i> , pages 8748–8763.	<i>Conference on Machine Learning</i> .	749
693	PmLR.	Hao Wu, Yuan Gao, Ruiqi Shu, Zean Han, Fan Xu,	750
694	Maziar Raissi, Paris Perdikaris, and George E Karni-	Zhihong Zhu, Qingsong Wen, Xian Wu, Kun Wang,	751
695	adakis. 2019. Physics-informed neural networks: A	and Xiaomeng Huang. 2025a. Turb-11: Achieving	752
696	deep learning framework for solving forward and	long-term turbulence tracing by tackling spectral bias.	753
697	inverse problems involving nonlinear partial differ-	<i>arXiv preprint arXiv:2505.19038</i> .	754
698	ential equations. <i>Journal of Computational physics</i> ,	Hao Wu, Yuan Gao, Ruiqi Shu, Kun Wang, Rui-	755
699	378:686–707.	jian Gou, Chuhan Wu, Xinliang Liu, Juncai He,	756
700	Bogdan Raonic, Roberto Molinaro, Tim De Ryck, To-	Shuhao Cao, Junfeng Fang, Xingjian Shi, Feng Tao,	757
701	bias Rohner, Francesca Bartolucci, Rima Alaifari,	Qi Song, Shengxuan Ji, Yanfei Xiang, Yuze Sun,	758
702	Siddhartha Mishra, and Emmanuel de Bezenac. 2023.	Jiahao Li, Fan Xu, Huanshuo Dong, and 7 others.	759
703	<i>Convolutional neural operators for robust and accu-</i>	2025b. Advanced long-term earth system forecast-	760
704	<i>rate learning of PDEs</i> . In <i>Thirty-seventh Conference</i>	by learning the small-scale nature. <i>arXiv preprint</i>	761
705	<i>on Neural Information Processing Systems</i> .	<i>arXiv:2505.19432</i> .	762
706	Stephan Rasp, Peter D Dueben, Sebastian Scher,	Hao Wu, Yuxuan Liang, Wei Xiong, Zhengyang Zhou,	763
707	Jonathan A Weyn, Soukayna Mouatadid, and	Wei Huang, Shilong Wang, and Kun Wang. 2024.	764
708	Nils Thuerey. 2020. Weatherbench: a bench-	Earthfarsser: Versatile spatio-temporal dynamical	765
709	mark data set for data-driven weather forecasting.	systems modeling in one model. In <i>Proceedings</i>	766
710	<i>Journal of Advances in Modeling Earth Systems</i> ,	<i>of the AAAI Conference on Artificial Intelligence</i> ,	767
711	12(11):e2020MS002203.	volume 38, pages 15906–15914.	768
712	Olaf Ronneberger, Philipp Fischer, and Thomas Brox.	Hao Wu, Shilong Wang, Yuxuan Liang, Zhengyang	769
713	2015. U-net: Convolutional networks for biomedical	Zhou, Wei Huang, Wei Xiong, and Kun Wang. 2023b.	770
714	image segmentation. In <i>Medical image computing</i>	Earthfarseeer: Versatile spatio-temporal dynamical	771
715	<i>and computer-assisted intervention–MICCAI 2015:</i>	systems modeling in one model. <i>AAAI2024</i> .	772
716	<i>18th international conference, Munich, Germany, Oc-</i>	Hao Wu, Wei Xion, Fan Xu, Xiao Luo, Chong	773
717	<i>ttober 5-9, 2015, proceedings, part III 18</i> , pages 234–	Chen, Xian-Sheng Hua, and Haixin Wang. 2023c.	774
718	241. Springer.	Pastnet: Introducing physical inductive biases for	775
719	Xingjian Shi, Zhoung Chen, Hao Wang, Dit-Yan Ye-	spatio-temporal video prediction. <i>arXiv preprint</i>	776
720	ung, Wai-Kin Wong, and Wang-chun Woo. 2015.	<i>arXiv:2305.11421</i> .	777
721	Convolutional lstm network: A machine learning ap-	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak	778
722	proach for precipitation nowcasting. <i>Advances in</i>	Shafraan, Karthik R Narasimhan, and Yuan Cao. 2022.	779
723	<i>neural information processing systems</i> , 28.		

- 780           React: Synergizing reasoning and acting in language  
781           models. In *The eleventh international conference on*  
782           *learning representations*.
- 783           Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing  
784           Sun, Tong Xu, and Enhong Chen. 2024. A survey on  
785           multimodal large language models. *National Science*  
786           *Review*, 11(12):nwae403.
- 787           Bing Yu, Haoteng Yin, and Zhanxing Zhu. 2018. Spatio-  
788           temporal graph convolutional networks: a deep learn-  
789           ing framework for traffic forecasting. pages 3634–  
790           3640.
- 791           Penghao Zhao, Hailin Zhang, Qinhan Yu, Zhen-  
792           gren Wang, Yunteng Geng, Fangcheng Fu, Ling  
793           Yang, Wentao Zhang, Jie Jiang, and Bin Cui. 2024.  
794           Retrieval-augmented generation for ai-generated con-  
795           tent: A survey. *arXiv preprint arXiv:2402.19473*.
- 796           Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang,  
797           Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen  
798           Zhang, Junjie Zhang, Zican Dong, and 1 others. 2023.  
799           A survey of large language models. *arXiv preprint*  
800           *arXiv:2303.18223*, 1(2).

## 801 **A Statistics for Data**

802 The evaluation of OMNIFLOW spans three mul-  
803 tiscale physical benchmarks, ranging from mi-  
804 croscopic fluid dynamics to global climate pat-  
805 terns. Table 3 summarizes the statistical details  
806 of these datasets. For 2D Turbulence, the model  
807 focuses on vorticity dynamics under continuous  
808 PDE constraints. The SEVIR dataset provides high-  
809 resolution regional convective patterns, while the  
810 ERA5 dataset serves as the global benchmark. For  
811 ERA5, we select a  $180 \times 360$  resolution and cu-  
812 rate a 200-day continuous sequence starting from  
813 January 1, 2020, covering 21 essential physical  
814 variables (including temperature, geopotential, and  
815 wind components across multiple pressure levels).

## 816 **B Experimental Setup and** 817 **Hyperparameters**

818 OMNIFLOW adopts a decoupled neuro-symbolic  
819 architecture as detailed in Table 4. Gemini 3 Flash  
820 serves as the cognitive core for ReAct planning,  
821 while the Neural Earth Simulator (NES) generates  
822 physical forecasts.

823 A critical aspect of our setup is the inference  
824 strategy: in 2D Turbulence, the model is trained on  
825  $1 \rightarrow 1$  step prediction but performs  $1 \rightarrow 99$  step  
826 long-term recursive extrapolation during testing  
827 to evaluate stability. For global forecasting, we  
828 employ a Perturbative Ensemble Strategy where  $K$   
829 members are generated via latent noise injection to  
830 quantify epistemic uncertainty ( $\sigma_{ens}$ ). The system  
831 also integrates a symbolic lens that automatically  
832 converts raw tensors into standardized units (e.g.,  
833 Kelvin to Celsius, Pascal to hPa). All experiments  
834 are averaged over 5 independent runs to ensure  
835 statistical significance.

## 836 **C Implementation Details for ERA5**

837 For the ERA5 benchmark, we extract 21 variables  
838 from the denormalized true labels and predictions.  
839 The variables include surface parameters (u10, v10,  
840 T2m, msl) and vertical profiles (U, V, T, Z, Q) at  
841 pressure levels ranging from 1000hPa to 100hPa.  
842 The simulation outputs are processed through a  
843 physical unit transformation layer to ensure align-  
844 ment with the knowledge retrieval loop’s expert  
845 reports. Each forecast frame is paired with a struc-  
846 tural JSON metadata file containing mean, max,  
847 min, and standard deviation for cross-modal rea-  
848 soning.

Table 3: Summary of Dataset Statistics (B6)

Dataset	Physical Regime	Resolution	Temporal Horizon	Key Variables	Evaluation Scale
2D Turbulence	Microscopic Flow	$128 \times 128$	100 timesteps	Vorticity	1,280 sequences
SEVIR	Regional Weather	$384 \times 384$	5-min intervals	VIL, IR, VIS, GLM	Standard Events
ERA5	Global Climate	$180 \times 360$	200-day forecast	21 variables (T, Z, U/V, Q)	Long-term Validation

Table 4: OMNIFLOW Configuration and Hyperparameters (C2)

Category	Component / Parameter	Specifications / Values
<b>Architecture</b>	Reasoning Core	Gemini 3 Flash (Agentic Brain)
	Numerical Simulator (NES)	Diffusion Transformer (DiT)
<b>Inference</b>	2D Turbulence Strategy	1 $\rightarrow$ 1 Training; 1 $\rightarrow$ 99 Recursive Testing
	ERA5 Strategy	200-day Autoregressive Forecast
	Ensemble Size ( $K$ )	8 – 32 (Adaptive)
<b>Hyperparameters</b>	Perturbation Factor ( $\lambda$ )	0.01 – 0.05
	Contrastive Temp ( $\tau$ )	0.07 (for Semantic-Symbolic Alignment)
	Independent Runs	5 (Mean $\pm$ Std)
<b>Control</b>	Consistency Check	Automated Unit Conversion ( $K \rightarrow ^\circ C$ , $Pa \rightarrow hPa$ )
	Evaluation Metrics	RMSE, SSIM, PSNR, Mech F1

## 849 D Structured Scientific Analysis Template

850 As illustrated in the below, OMNIFLOW generates  
851 a structured, multi-faceted scientific report that  
852 bridges the gap between raw numerical tensors and  
853 interpretable expert analysis. The report template  
854 is designed to provide human-auditable reasoning  
855 chains, ensuring physical consistency and opera-  
856 tional utility. The template consists of four primary  
857 modules:

- 858 • **Executive Summary:** This section provides  
859 a high-level synoptic overview of the phys-  
860 ical state (e.g., global atmospheric circula-  
861 tion). It contextualizes the temporal and spa-  
862 tial bounds of the data and identifies dominant  
863 physical phenomena, such as significant pres-  
864 sure gradients or hemispheric dynamics.
- 865 • **Statistical Overview (Data Analysis):** To  
866 maintain scientific rigor, this module extracts  
867 precise quantitative metrics from the predicted  
868 flow fields. Key indicators include the *Global*  
869 *Mean*, *Extreme Minimum/Maximum* (with 4-  
870 decimal precision), and *Variability (Standard*  
871 *Deviation)*. This ensures that the agent’s rea-  
872 soning is grounded in exact numerical evi-  
873 dence rather than vague qualitative assess-  
874 ments.
- 875 • **Spatial Pattern Analysis (Visual Inter-**  
876 **pretation):** Leveraging the *Symbolic Lens*,  
877 the agent translates visual features (e.g.,  
878 heatmaps) into topological linguistic descrip-  
879 tors. It identifies high-pressure zones (bright  
880 regions) and low-pressure vortices (dark re-  
881 gions), mapping them to known physical struc-  
882 tures like continental air masses or subtropical  
883 ridges.
- 884 • **Meteorological Insights & Conclusion:** The  
885 final module synthesizes the quantitative and  
886 spatial data to derive expert-level conclusions.  
887 It assesses the synoptic situation, evaluates  
888 potential impacts (e.g., geostrophic wind in-  
889 tensity or winter storm activity), and provides  
890 a basis for downstream decision-making and  
891 risk assessment.

## 1. Executive Summary

On January 1, 2020, the global atmospheric circulation presented a robust mean sea level pressure field. The analysis highlights significant pressure gradients, particularly characteristic of Northern Hemisphere winter dynamics, with deep low-pressure systems driving high variability.

## 2. Statistical Overview (Data Analysis)

- **Global Mean (1008.6248779296875 hPa):** Near the standard atmospheric average of 1013 hPa, indicating a balanced global mass distribution.
- **Extreme Minimum (953.4845581054688 hPa):** This extremely low value strongly suggests the presence of a deep extratropical cyclone or a polar vortex fragment, likely located in the high latitudes (North Atlantic or Southern Ocean).
- **Extreme Maximum (1045.9444580078125 hPa):** Indicates a powerful anticyclone (high-pressure system), likely associated with cold continental air masses (e.g., Siberian High) or subtropical ridges.
- **Variability (Std 12.760520935058594 hPa):** A standard deviation of  $\sim 12.8$  hPa reflects an active weather pattern with steep pressure gradients, typical for January.

## 3. Spatial Pattern Analysis (Visual Interpretation)

Observing the heatmap:

- **High Pressure Zones (Yellow/Bright):** The brightest regions likely correspond to the subtropical high-pressure belts and the continental highs over Eurasia (Siberian High). These areas represent stable, sinking air.
- **Low Pressure Zones (Blue/Dark):** The dark vortices visible in the mid-to-high latitudes confirm the statistical minimums. These “swirls” are visually consistent with storm tracks in the North Atlantic and Pacific, as well as the circumpolar trough around Antarctica.

## 4. Meteorological Insights & Conclusion

The ERA5 Ground Truth data confirms a highly active synoptic situation. The span of nearly 92 hPa between the minimum and maximum pressure drives strong geostrophic winds. The data implies intense winter storm activity in the Northern Hemisphere, while the tropics likely remain dominated by standard easterly trade wind gradients.