



Can We Trust LLMs for Complex Earth System Model Analysis?

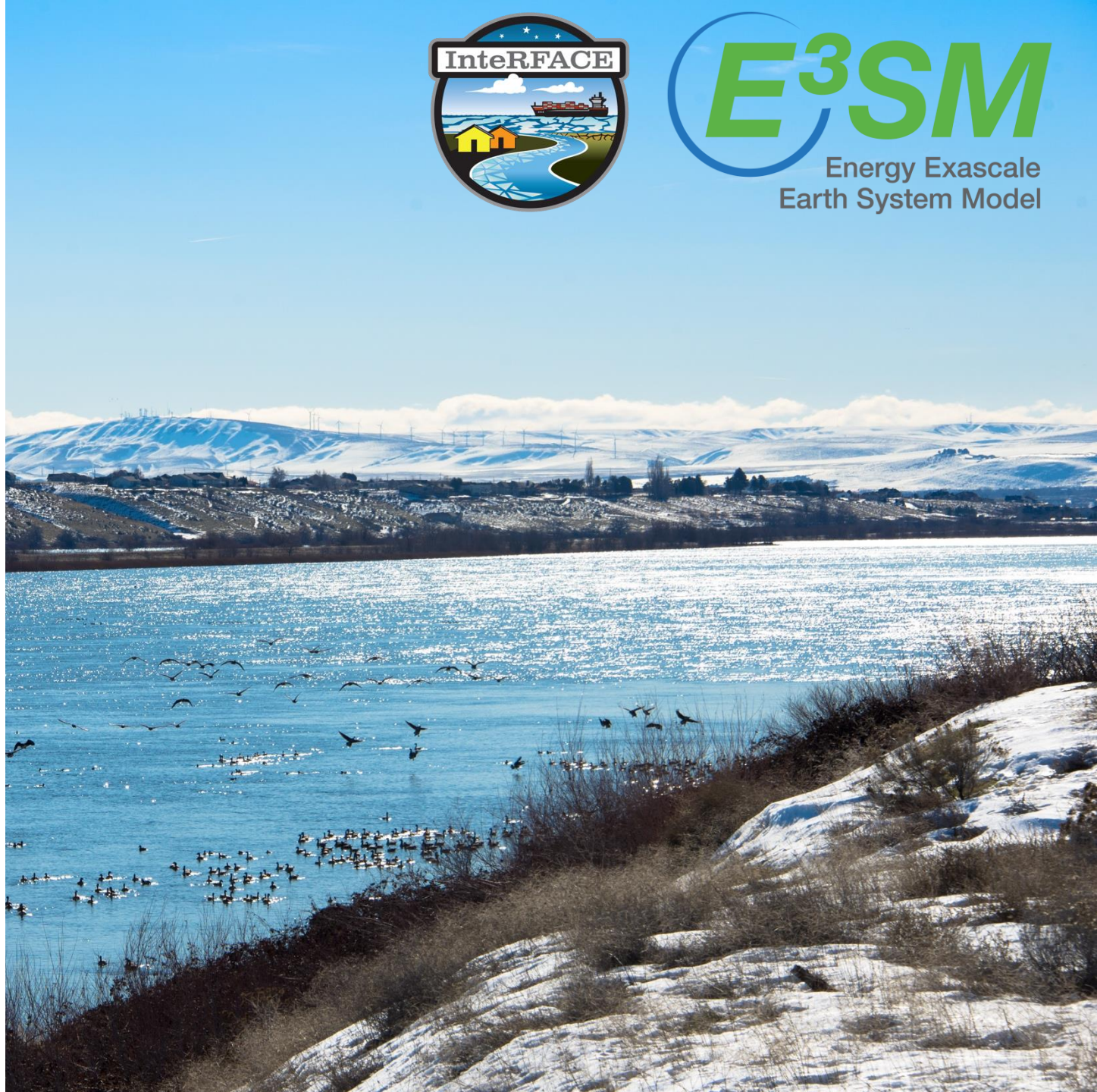
Tian Zhou

HydroML26

5/19/2026



PNNL is operated by Battelle for the U.S. Department of Energy



**LLMs can now write
working scientific scripts.**

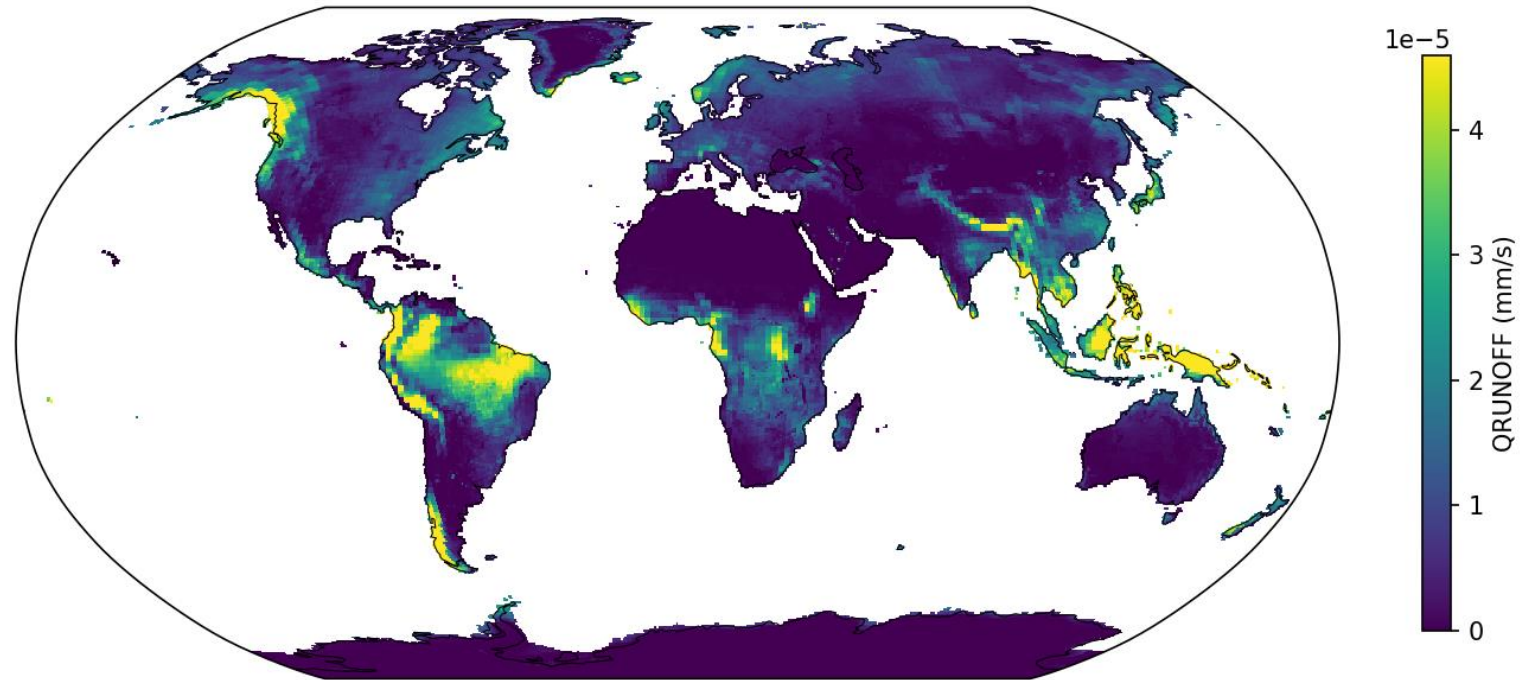
**LLMs can now write
working scientific scripts.**

**But is the method
scientifically valid?**

You are a scientific programmer analyzing ESM output. Please write a complete, self-contained Python script that performs the requested analysis.

Task: Using E3SM land model monthly output, extract the climatological mean runoff over 1985–1989. The ELM data path is XXX. Extract the global QRUNOFF field for 1985–1989 and compute the area-weighted global mean. Produce a map of the mean runoff field with the global statistics overlaid.

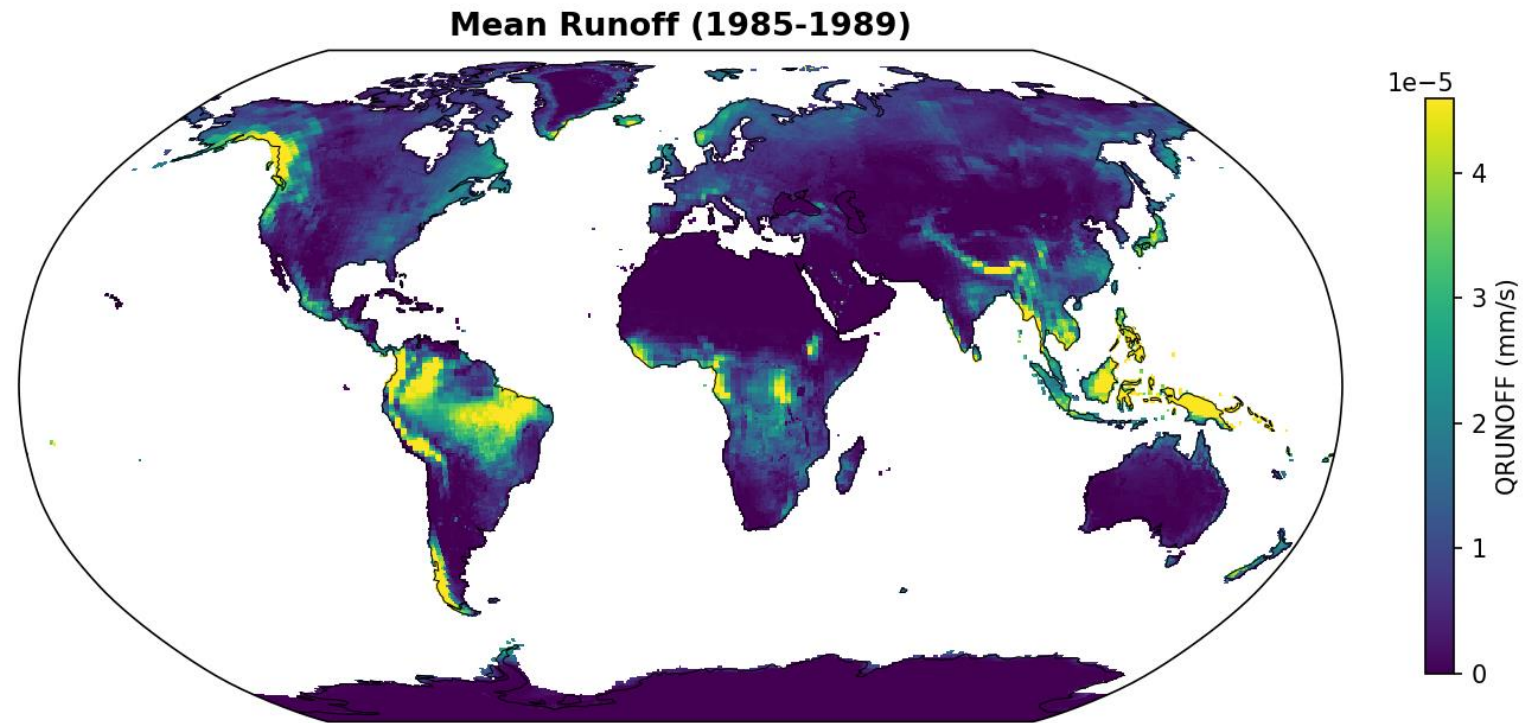
Mean Runoff (1985-1989)



global: $6.75e-06$ mm/s

You are a scientific programmer analyzing ESM output. Please write a complete, self-contained Python script that performs the requested analysis.

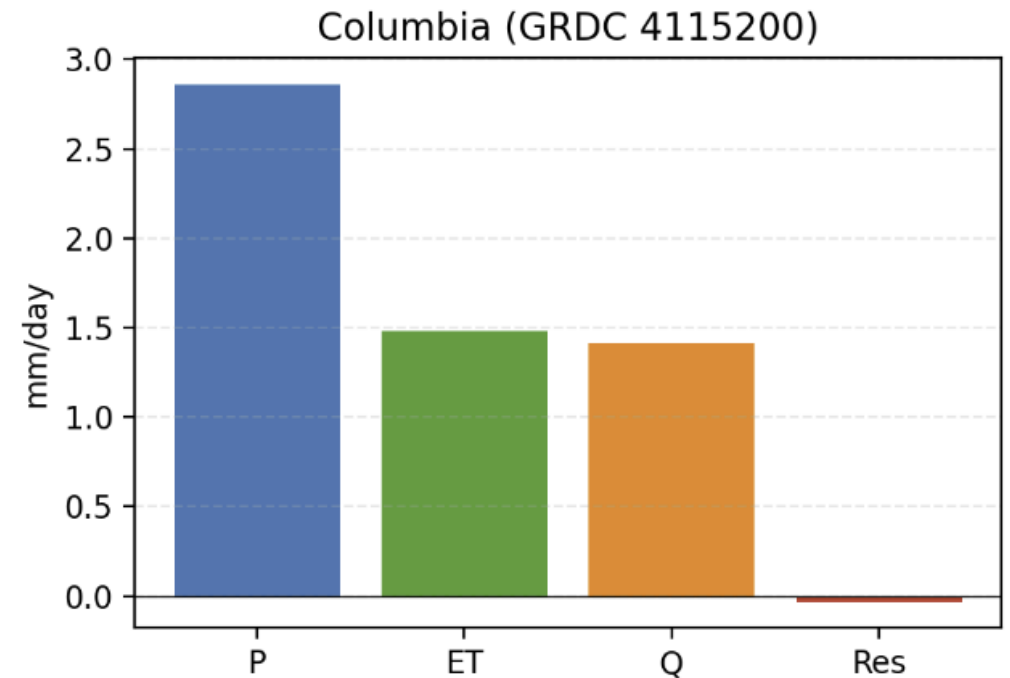
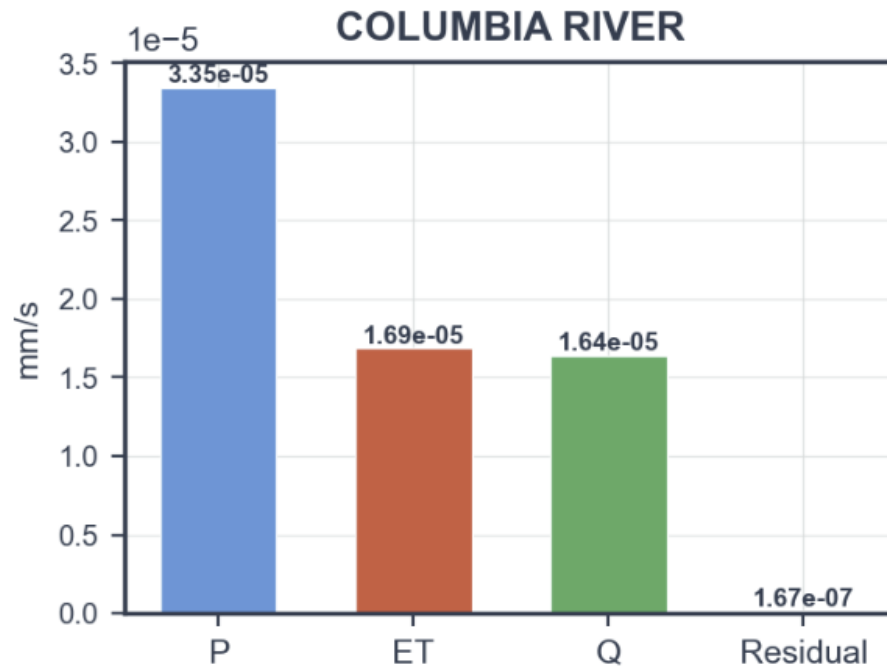
Task: Using E3SM land model monthly output, extract the climatological mean runoff over 1985–1989. The ELM data path is XXX. Extract the global QRUNOFF field for 1985–1989 and compute the area-weighted global mean. Produce a map of the mean runoff field with the global statistics overlaid.



Silent Failure!

global: 6.75e-06 mm/s

Silent Failure Could Lead to Opposite Conclusions



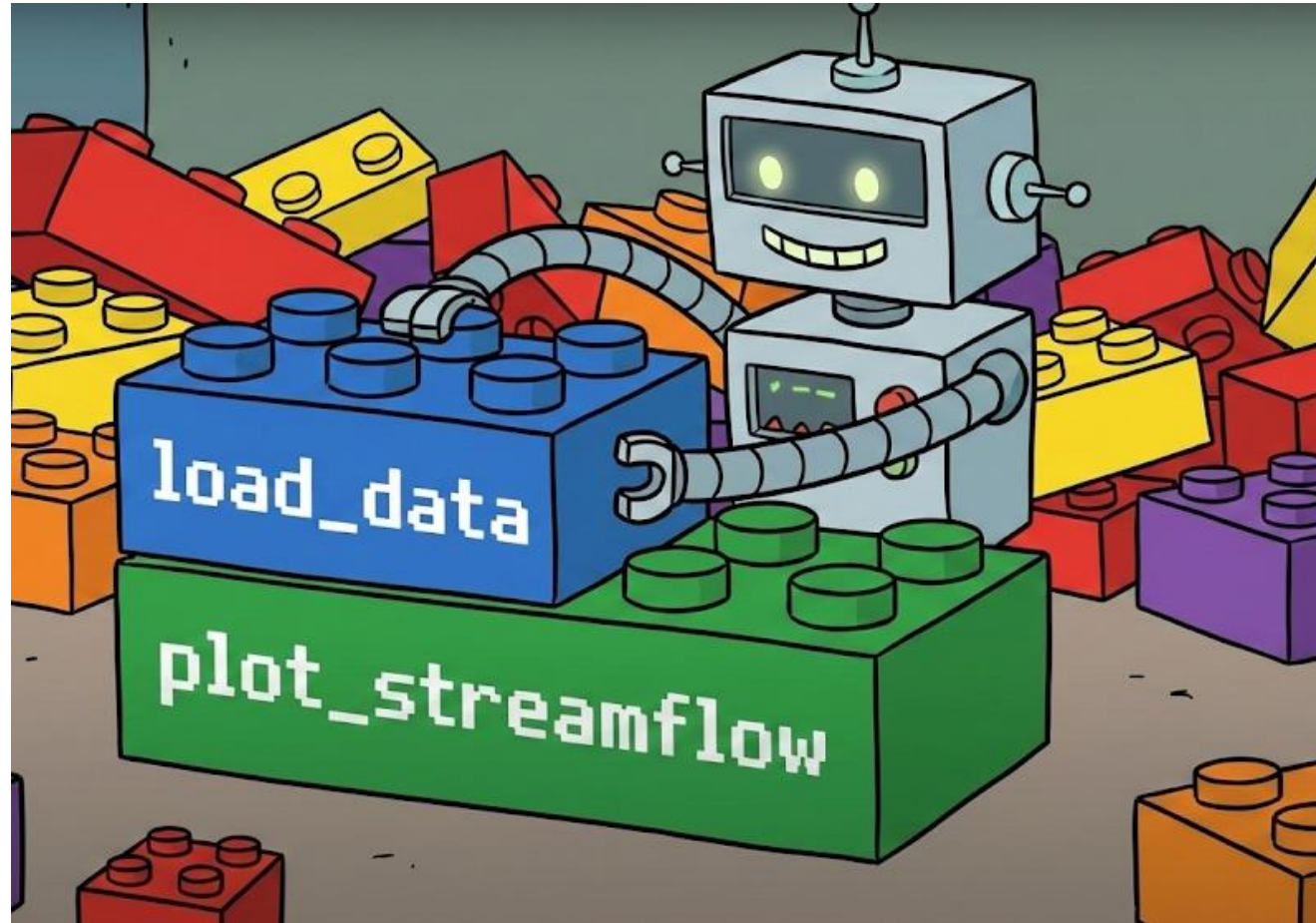
The auto-generated script omits land-fraction weighting, leading to negative water budget residual calculation.

Silent failures are runs that
succeed at the software level but
fail at the scientific level.

Silent failures are runs that succeed at the software level but fail at the scientific level.

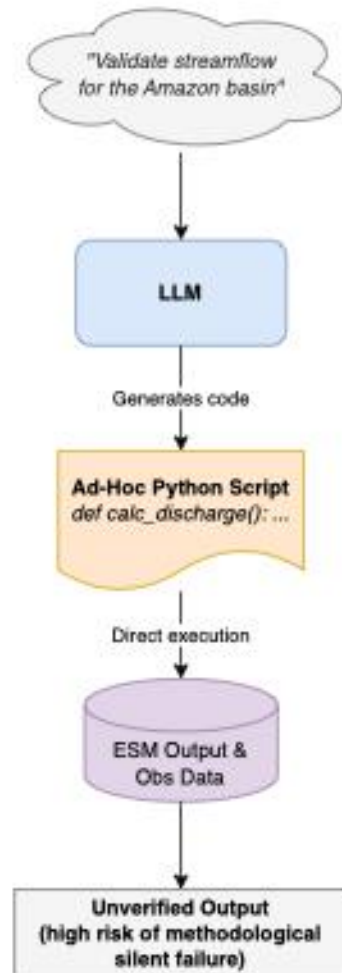
What should we let LLM to decide?

Let the LLM compose, not code

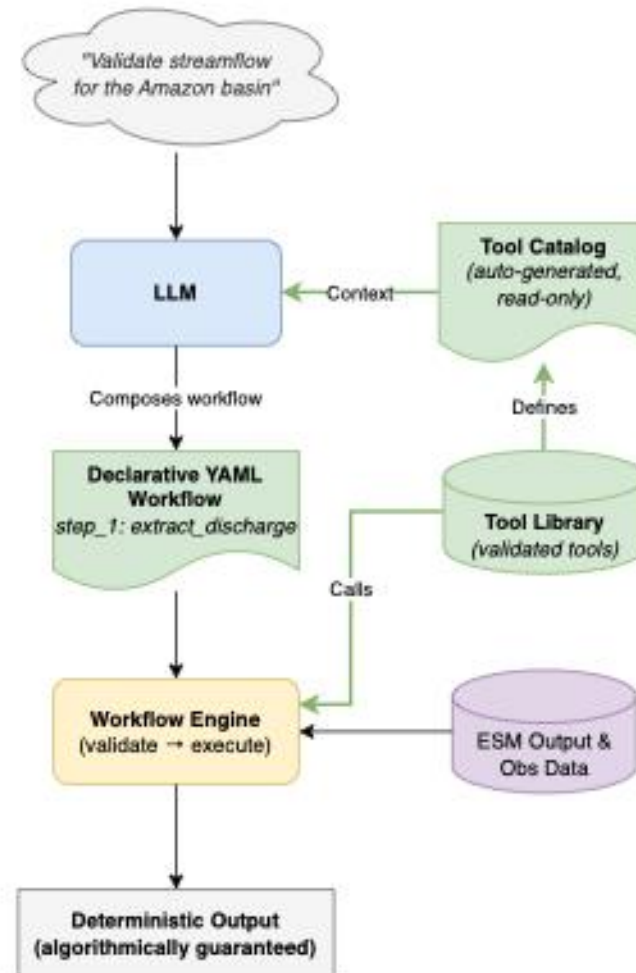


ESFlow: a Module-Grounded Workflow

(a) Unconstrained Code Generation



(b) Module-Grounded Workflow Composition



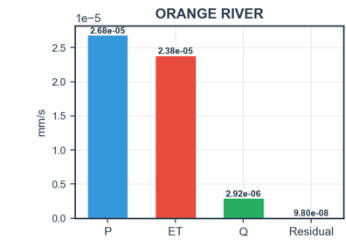
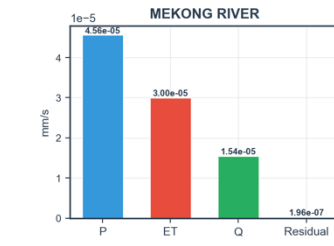
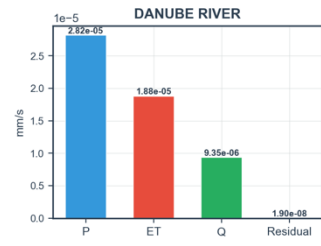
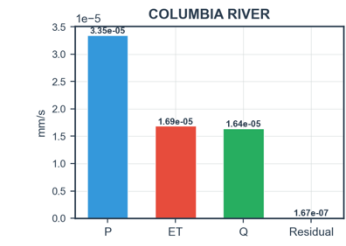
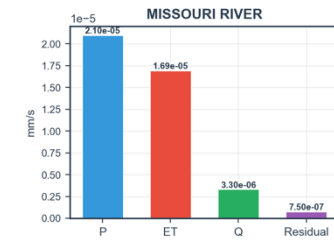
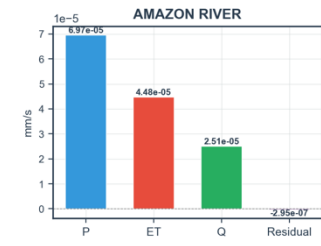
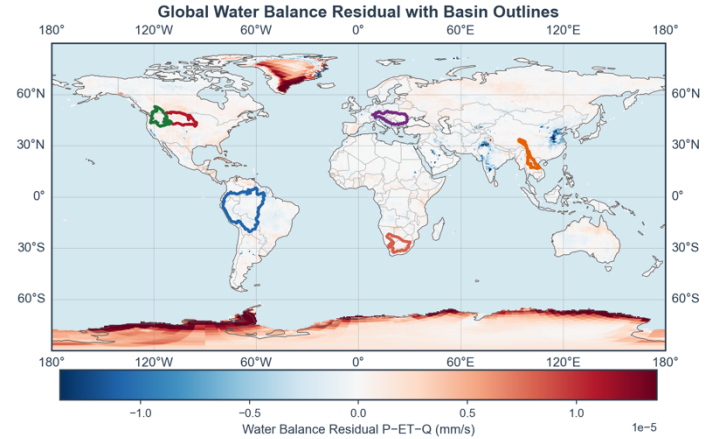
- Validated Python tools (LLMs never see)
- Tool Library describes tool function and I/O (LLM context)
- LLM generates workflow in YAML format
- A `run_workflow` engine executes the workflow by calling tools

Benchmark tests

- Two test types: baseline (direct Python generation) & module-grounded (ESFlow framework).
- Six LLMs evaluated: Phi-4, Haiku 4.5, o4-mini, Gemini 2.5 Flash, GPT-5, Opus 4.6.
- Seven hydrology diagnostic tasks – ranging from simple to complex.
- Four repeats per test run.

Benchmark tests (Task 6)

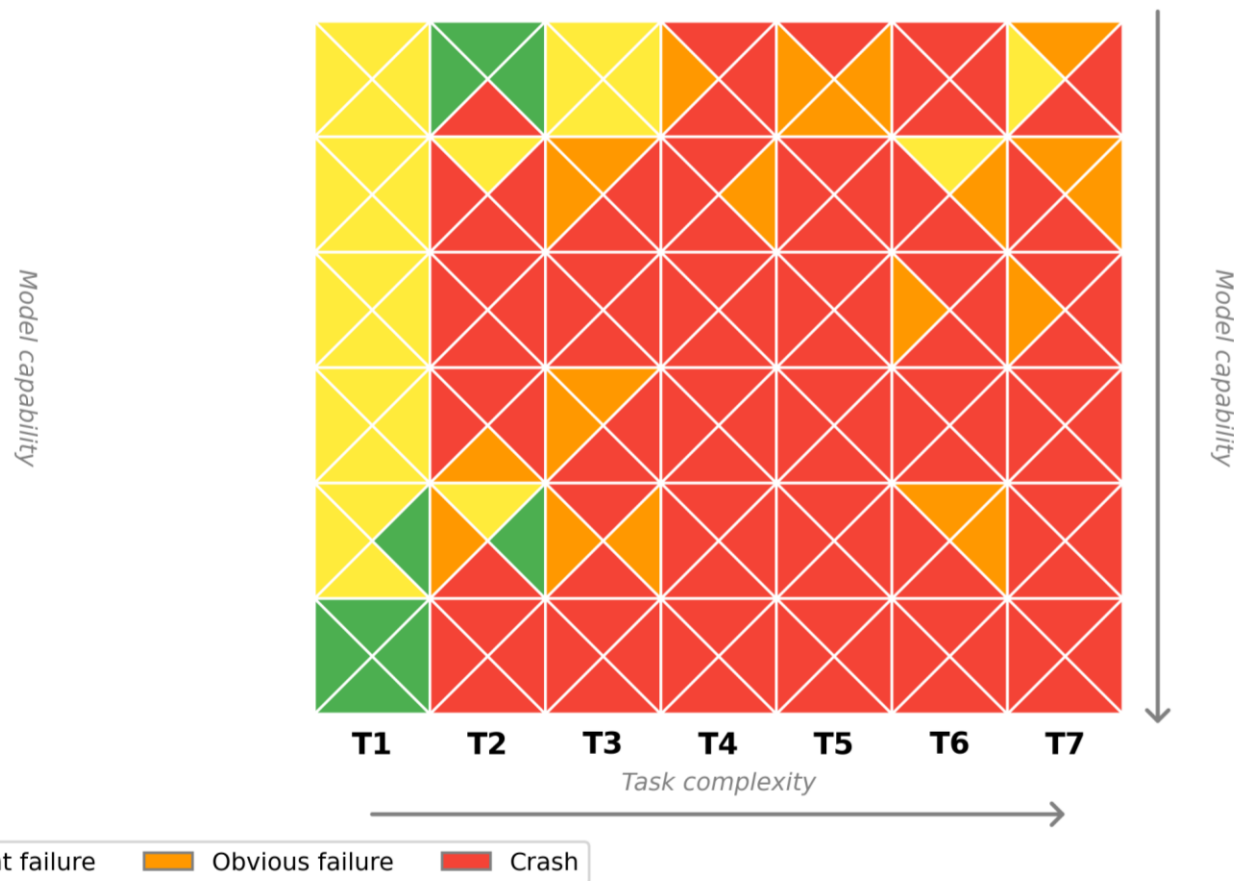
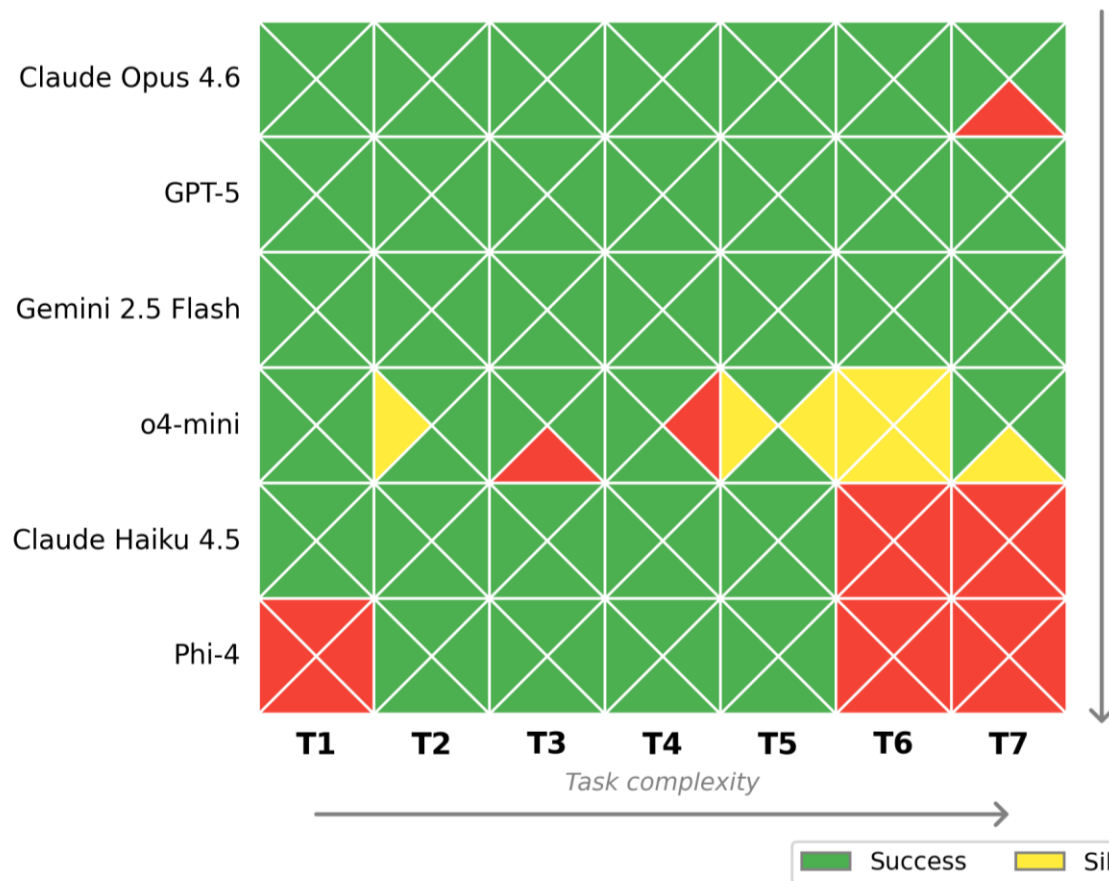
T6 Reference: Water Balance



Single shot results

(a) Module-grounded (YAML Workflow)

(b) Baseline (Python Code-Gen)



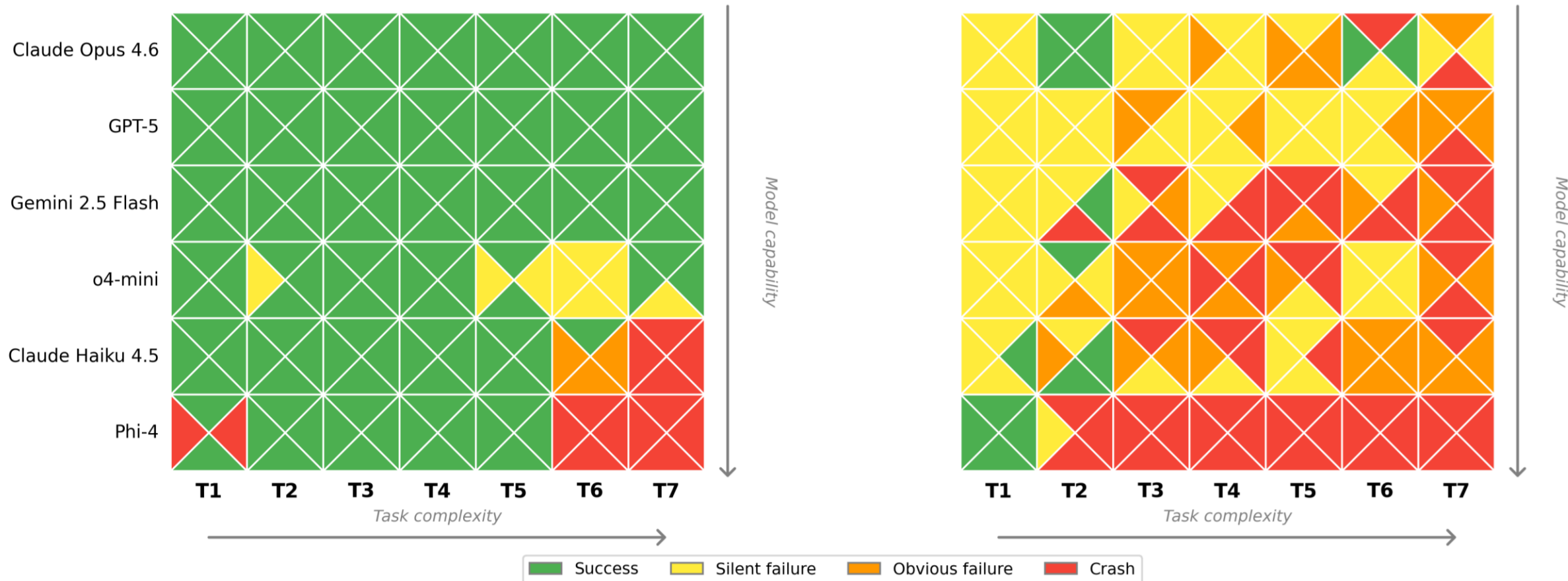
82% Success, 5% Silent Failure

5% Success, 16% Silent Failure

Self debug (up to 3 rounds) results

(a) Module-grounded (YAML Workflow)

(b) Baseline (Python Code-Gen)



85% Success, 5% Silent Failure

9% Success, 40% Silent Failure

Can we trust LLMs for complex ESM analysis?

- **Unconstrained code generation is unreliable for ESM analysis.**
- **Module-grounded composition achieves high-fidelity results.**
- **Self-debugging is safe only when the LLM is constrained.**



Thanks for your
attention!



Can We Trust LLMs for Complex Earth System Model Analysis? Silent Failure and Evidence from Module-Grounded Benchmarking

Tian Zhou¹, Yun Qian¹, and L. Ruby Leung¹

¹Pacific Northwest National Laboratory, Richland, WA, USA

Correspondence: Tian Zhou (tian.zhou@pnnl.gov)

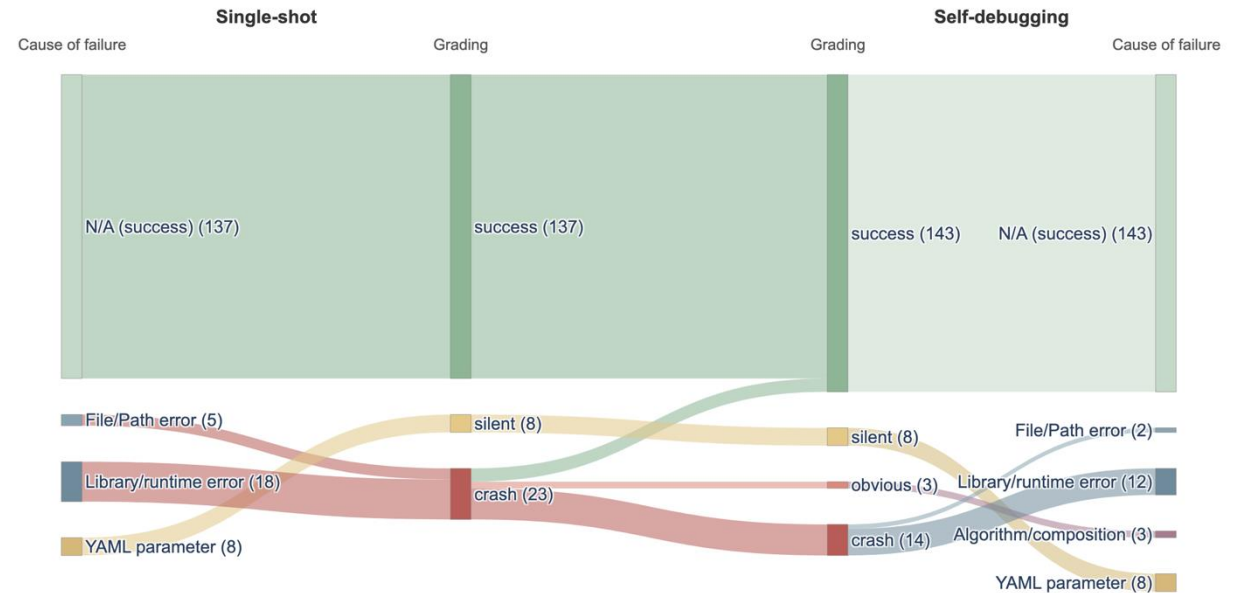
Abstract.

Large language models (LLMs) are becoming increasingly capable of complex scientific scripting, but this growing robustness creates a paradox: the more trustworthy their outputs appear, the more easily scientifically incorrect results can pass unnoticed. In Earth system model (ESM) analysis, such silent failures are more dangerous than visible crashes because they produce plausible figures and statistics that may be accepted without detailed inspection. We address this risk with ESFlow, a module-grounded agentic AI framework that constrains the LLM to compose workflows from validated analysis tools rather than generate arbitrary code. The LLM reads an auto-generated, self-describing catalog and outputs a YAML (human-readable data-serialization) workflow, which is then executed by a deterministic engine. We demonstrate this framework with a validated tool library for Energy Exascale Earth System Model (E3SM) land surface hydrology diagnostics in a benchmark spanning seven analysis tasks and six contemporary LLMs. Across both single-attempt runs and runs augmented with automatic self-debugging, the module-grounded approach attains an overall success rate above 80 %, maintains a low and stable silent-failure rate, and reaches 100 % success for the three high-capability models, whereas unconstrained Python code generation succeeds in only about 5 % of runs and sees its silent-failure rate rise from roughly 16 % to about 40 % under self-debugging. These results suggest that increasing LLM capability does not remove the reliability problem in scientific scripting; it makes silent failures more consequential by making incorrect outputs more convincing. The answer to the trust question posed in the title is therefore conditional: unconstrained code generation is not trustworthy for complex ESM analysis, whereas module-grounded workflow composition can be highly reliable for frontier models and remains substantially more robust under iterative self-debugging. By shifting the LLM's role from code generation to the composition of trusted tools, this framework provides a safer, more scalable architecture for AI-assisted scientific discovery that is aligned with FAIR (findable, accessible, interoperable, and reusable) principles.

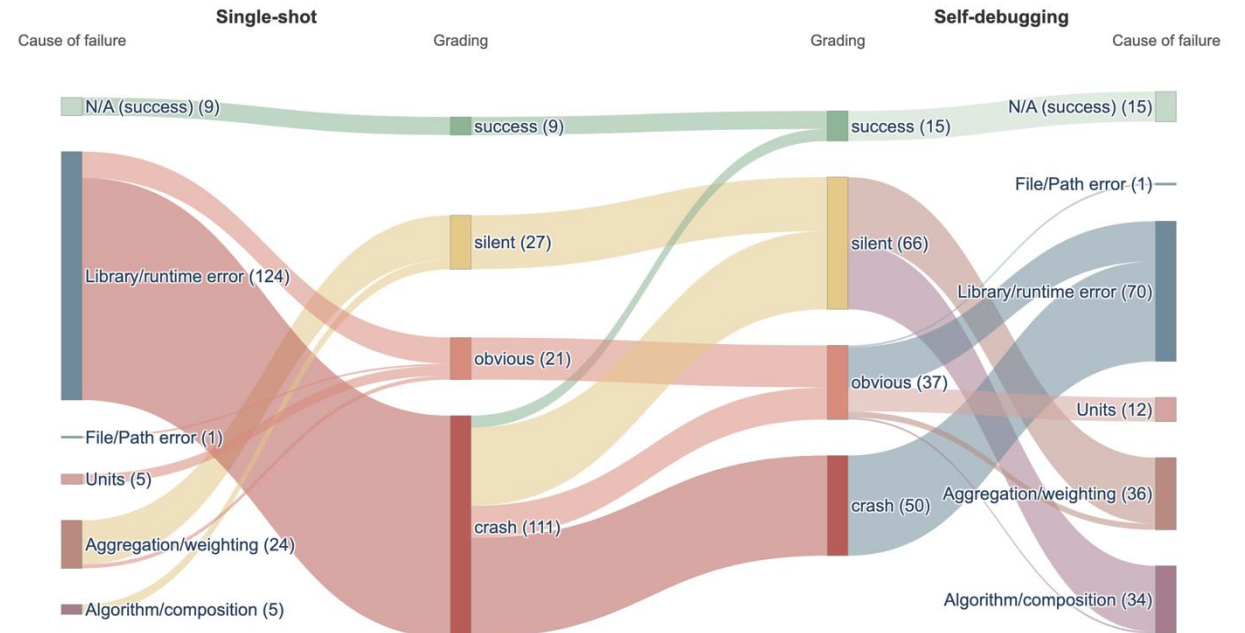
Backup Slide

Flow of non-success runs through single-shot grading and post-self-debug

(a) Module-grounded (YAML Workflow)



(b) Baseline (Python Code-Gen)



Backup Slide

Root cause of the non-success runs

Category	Description and representative examples
File/Path error <i>(execution-level)</i>	Crash from a wrong metadata path, missing input file, or unsupported output location. <ul style="list-style-type: none"> • Wrong gauge-file path (Phi-4, T1) • Output organised into subdirectories the engine does not create (Opus, T7) • Incorrect observation-metadata path (Haiku, T4)
Library/runtime error <i>(execution-level)</i>	Language-level or library-level failure during script or tool execution. <ul style="list-style-type: none"> • Calendar-type mismatch when coercing E3SM no-leap dates (multiple models, T4/T5) • NumPy shape mismatch during gauge-to-grid matching, e.g. (360,) vs. (720,) (GPT-5, T4) • Lazy-array evaluation failure on scalar extraction (Gemini, T7) • Polygon-clipping error on self-intersecting basin geometries (multiple models, T7) • Missing or incorrect import statements (Phi-4, T6)
Units <i>(methodological)</i>	Missing or partial unit conversion that produces plausible but numerically wrong output. <ul style="list-style-type: none"> • Converts some variables to mm day^{-1} but leaves others in mm s^{-1} (multiple models, T7) • Omits MODIS $\text{kg m}^{-2} \text{s}^{-1}$ to mm day^{-1} conversion (multiple models, T3) • Applies wrong conversion factor ($\times 86.4$ instead of $\times 86400$) (o4-mini, T7)
Aggregation / weighting <i>(methodological)</i>	Incorrect temporal aggregation, spatial weighting, or land/ocean masking. <ul style="list-style-type: none"> • Daily-to-monthly aggregation inconsistent between model and observations (multiple models, T2/T4) • Missing land-fraction weighting, causing 1–2 % drift at coastal cells (multiple models, T6) • Area-weighted mean over all cells including ocean zeros, diluting land-only values (multiple models, T3) • Drops cells with negative runoff before averaging, biasing upward (Gemini, T6)
Algorithm / composition <i>(methodological)</i>	Wrong analytical method, missing workflow step, or incorrect metric definition. <ul style="list-style-type: none"> • Bilinear interpolation instead of conservative remapping (multiple models, T3) • Un-normalised Wasserstein distances reported as raw values, e.g. 54 144 vs. reference 0.32 (Opus, T7) • P–ET reported instead of P–ET–Q, omitting the runoff-subtraction step (Haiku, T6) • Hand-rolled spherical area weights instead of model-native area variable (GPT-5, T6) • Longitude convention mismatch ($0\text{--}360^\circ$ vs. $\text{--}180\text{--}180^\circ$) producing all-NaN fields (Phi-4, T3)
YAML parameter <i>(parameter; module-grounded only)</i>	Structurally valid workflow that invokes a validated tool with incorrect arguments; the tool executes normally but on wrong inputs. <ul style="list-style-type: none"> • <code>years: [1985, 1989]</code> (two-element list) instead of the full five-year range, loading only 2 of 5 years (o4-mini, T5/T6/T7) • Wrong variable name passed to a correctly selected extraction tool (o4-mini, T4)