

✓-February Flow

Data Components:

Code:

The-stack-v2

CodeText:

SE, whatever we're scraped

WebText:

HQ DCLM

DATA MIXES

~85%	Source Code	DeepSeek Coder
~10%	CodeText	
~5%	Webtext	

~85%	The-stack-v2	StarCoder 2
~15%	CodeText	
~0%	Webtext	

~100%	Source Code	Arctic
-------	-------------	--------

P1: 100% Source Code] Granite
P2: 80% Code
20% language]

Code Data Recipe [Stacked]

- 1) Order by Repo ✓
- 2) Call Heuristic Filters X
- 3) Group by Repo, lang → minhash ✓
- 4) Pack into Repo-level docs △
- 5) Select PL's △
- ⋮
- 6) Pack into FIM tokens* X

✓: Eng Done

*not critical

X: Eng definitely NOT done

△: So So easy

Use Preprocessed

cootext, webtext

ARCH + TRAINING

- Pick Arch like OLMo - IB
~OR~ replicate a 3D model
- Follow standard LR flow

EVAL:

Hacky nonsense for now