

Debugging Guide & Best Practices



"If you don't audit the gradients, you are training on hope."

This document outlines the standard fixes for common pathologies in [PROJECT_TYPE] models.

1. The Physiology of Training (Weights & Gradients)

A. The Weight Update Ratio

We must ensure weights are actually moving, but not too fast.

$$\text{Ratio} = \frac{\|\eta \cdot \nabla W\|}{\|W\|}$$

- **Target: 1e-3** (The "Goldilocks Zone").
- **< 1e-5**: Model is frozen (Check Learning Rate, Detached Graphs).
- **> 1e-1**: Model is unstable (Check Clipping, Batch Size).

B. Gradient Norm Topology

- **Norm = 0.0: Broken Graph.** Check `requires_grad`, `.detach()`, or index operations.
- **Norm > 10.0: Exploding.** Use `clip_grad_norm_`.
- **Norm = NaN: Overflow.** Check for division by zero or `log(0)`.

2. Shape Hygiene



Symptom: `RuntimeError: The size of tensor a (X) must match the size of tensor b (Y)`.

Best Practices

1. Assert Inputs:

```
B, T, D = x.shape  
assert D == self.hidden_dim, f"Mismatch: {D} != {self.hidden_dim}"
```

2. Use Einops: Avoid `x.view()`. Use `rearrange(x, 'b t (h d) -> ...')` to be explicit.

3. Convergence Pathology



1. The "Hockey Stick": Loss drops instantly then flatlines.

- *Diagnosis:* LR too high. Local Fast Basin.

2. The "Slow Burn": Loss decreases linearly.

- *Diagnosis:* LR too low.

3. The "NaN Death": Loss becomes NaN at step N.

- **Diagnosis:** Optimization instability or Normalization Drift.

4. GPU/Hardware Tips

- **OOM (Out of Memory):** Reduce Batch Size, use Gradient Accumulation.
- **Slow Training:** Check CPU-GPU sync points (e.g. `print(tensor)` inside loop).

5. Agent Debugging Workflow

When you (an AI agent) encounter training issues, follow this decision tree:

Step 1: Gather Evidence (`view_file`, `run_command`)

```
# Check last 50 lines of log
tail -n 50 logs/train.log

# Check GPU memory
nvidia-smi

# Check if process is still running
ps aux | grep train
```

Step 2: Classify the Failure Mode

Symptom	Tools to Call	What to Check
Loss = NaN	<code>view_code_item</code> : Normalization layers, loss functions	Look for <code>log(0)</code> , division by zero, <code>sqrt(negative)</code>
Loss doesn't decrease	<code>view_code_item</code> : Optimizer init, learning rate schedule	Verify <code>requires_grad=True</code> , check LR value
OOM Error	<code>view_file</code> : Model config	Reduce <code>batch_size</code> or <code>hidden_dim</code> , use gradient checkpointing
Shape Mismatch	<code>view_code_item</code> : Forward pass where error occurs	Add <code>assert</code> statements for tensor shapes
Slow throughput	<code>grep_search</code> : Search for <code>.item()</code> or <code>print(tensor)</code> in training loop	Remove CPU-GPU sync points
MCP Server Not Found	Use <code>view_file</code> or <code>view_file_outline</code>	You tried to use <code>read_resource</code> or <code>list_resources</code> on a local path. Filesystem paths are not MCP servers.
Artifact Access Error	Use absolute path to <code>.gemini/antigravity/...</code>	Artifacts are local files. Access them via <code>view_file</code> .

Step 3: Implement Fix (`replace_file_content`)

Example: Loss = NaN at step 500

1. **view_file**: Check log for exact error line
2. **view_code_item**: View the layer where NaN occurs
3. **Hypothesis**: Entropy calculation has $\log(0)$
4. **replace_file_content**: Add epsilon: `entropy = -torch.sum(p * torch.log(p + 1e-8))`
5. **run_command**: Relaunch training from checkpoint

Step 4: Document in Journal (**replace_file_content**)

Add to **journal.md**:

```
## Debugging Log

### [DATE] - NaN Loss at Step 500
**Symptom**: Loss became NaN during entropy calculation
**Root Cause**: Missing epsilon in log operation
**Fix**: Added `1e-8` epsilon to prevent `log(0)`
**Resolution**: Training resumed successfully from checkpoint
```

Step 5: Create Smoke Test (**write_to_file**)

Prevent regression:

```
# tests/test_entropy_stability.py
def test_entropy_no_nan():
    p = torch.tensor([0.0, 1.0]) # Edge case
    entropy = -torch.sum(p * torch.log(p + 1e-8))
    assert not torch.isnan(entropy), "Entropy should not be NaN"
```

Remember: Every bug you fix should generate:

1. A journal entry (documentation)
2. A test case (prevention)