# Toward Reproducible Cross-Backend Compatibility for Deep Learning: A Configuration-First Framework with Three-Tier Verification

Z. Li
Department of Computer Science
Dalhousie University
Email: zehua.li@dal.ca

*Abstract*—**Cross-backend behavioral drift threatens the reproducibility of deep learning systems deployed on CPU, GPU, and compiled runtimes. We study three questions: (RQ1) How large is cross-backend behavioral drift under practical tolerances? (RQ2) Which models/tasks are most prone to cross-backend inconsistencies? (RQ3) Where (which layers) does divergence first emerge? We propose a configuration-first framework that decouples experiments from code via YAML, supports both library and repository models, and verifies outputs with a three-tier strategy: tensor closeness, activation alignment, and task-level metrics. The framework emits structured JSONL logs and integrates into CI. Given the compute constraints of this exploratory study, we emphasize end-to-end and task-level agreement while retaining activation-level probing as an optional capability. Across four tolerance settings (atol $\in \{10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}\}$) and $672$ cross-backend checks, we observe $484$ passes ($72.0\%$) in aggregate, with most discrepancies concentrated at tighter tolerances. To our knowledge, this is the first unified framework that systematically *quantifies and mitigates* cross-backend drift via a configuration-first, three-tier protocol.**

*Index Terms*—**Reproducibility, cross-backend drift, numerical stability, deep learning systems, PyTorch, deterministic adapters.**

Fig. 1: Pipeline: YAML $\rightarrow$ Loader $\rightarrow$ Preprocess $\rightarrow$ Exec (ref,tgt) $\rightarrow$ Verify $\rightarrow$ Reports.

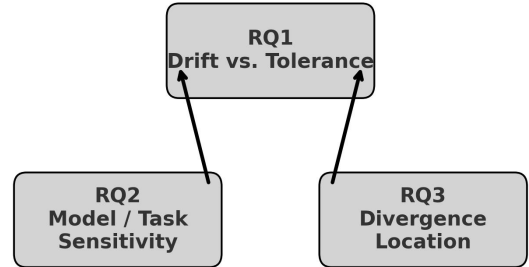We also provide an at-a-glance summary of our evaluation axes:



Fig. 2: Overview of RQ1–RQ3: tolerance sweep (RQ1), model$\times$backend sensitivity (RQ2), and divergence localization (RQ3).

## I. INTRODUCTION

### A. Motivation and Stakes

Cross-backend discrepancies can arise from kernel implementations, precision modes, autotuning, and graph rewrites. Even minor numerical perturbations can alter post-processing outcomes in detection or segmentation. Such changes may be consequential in safety-critical deployments (e.g. autonomous driving, medical imaging). For example, in a detection pipeline, CPU vs. GPU/compiled runs produced different pre-NMS orderings of candidate boxes, changing final picks despite tensor-level differences within $10^{-5}$. Enforcing a deterministic sort prior to NMS eliminated this inconsistency.

### B. Research Questions and Challenges

We address:

- **RQ1:** How large is cross-backend behavioral drift under practical tolerances?
- **RQ2:** Which models/tasks are most prone to cross-backend inconsistencies?
- **RQ3:** Where (which layers) does divergence first emerge?

Key challenges include nondeterminism, operator coverage gaps, and downstream post-processing variability, each of which can amplify small numerical differences into task-level failures.

### C. Contributions

- **Configuration-first methodology.** A backend-agnostic runner that decouples experiment design from implementation via YAML, improving reproducibility and reuse.
- **Three-tier verification protocol.** A unified evaluation at tensor/activation/task levels with deterministic adapters for post-processing.
- **Empirical characterization.** A study across models $\times$ backends $\times$ tolerances, together with a succinct failure taxonomy and latency analysis.
- **Actionable mitigations.** Deterministic pre-NMS sorting, selective eager fallbacks, and FP32 enforcement that significantly improve agreement with minimal overhead.

*Claim.* To our knowledge, this is the first unified protocol that directly links tensor-level drift to task-level outcomes and validates fixes under a common configuration-first design.

## II. Related Work

Prior work on testing neural networks has largely focused on input-space robustness or interface-level checks (DeepXplore, DeepTest, TensorFuzz, Mist), whereas compiler efforts (TVM, XLA, Glow, Inductor) emphasize transformation soundness. Deterministic flags, seed control, and activation probing are common practices; our contribution is to systematize these within a unified protocol aimed at cross-backend agreement. Unlike input fuzzing [1], [2], [3], [4] and compiler validation [5], our approach explicitly links tensor-level drift to task-level outcomes via deterministic adapters, aligning with reproducibility guidelines [6] and addressing deployment inconsistencies in heterogeneous runtime settings.

## III. Problem Formulation and Methodology

### A. Compatibility Criterion

Let $M$ be a model with fixed weights, $B = \{b_1, \dots, b_k\}$ a set of backends, and $y_i = f(M, x; b_i)$ the corresponding outputs. We declare tensor-level compatibility if

$$\|y_i - y_j\|_\infty \leq \text{atol} + \text{rtol} \cdot \|y_i\|_\infty. \tag{1}$$

We additionally track MAE, $p95$ error, and task metrics (Top-1/Top-5, mAP, mIoU) to avoid false alarms from benign permutations and to better reflect end-task fidelity.

### B. Three-Tier Verification

**Tier-1 (Tensor):** Eq. (1) and error statistics.
**Tier-2 (Activation):** lightweight hooks for layerwise probing to localize the earliest divergence.
**Tier-3 (Task):** deterministic post-processing adapters (e.g. sorting keys before NMS) and metric-level agreement.

**Scope note (RQ3).** We retain activation-level instrumentation, and employ it selectively to demonstrate feasibility. A comprehensive activation survey is deferred given the cost of large-scale multi-backend sweeps.

### C. Configuration-First Execution

Experiments are YAML-driven: model source (`library`/`repo`), preprocessing (means/std, resize), backends/compile options, and tolerances. This design decouples experiment specification from code, facilitating replication and extension.

---

**Algorithm 1** Compatibility Runner (Sketch)

---

1: **Input:** YAML configs $\mathcal{C}$, backends $B$, tolerances $(\text{atol}, \text{rtol})$
2: **for** config $c \in \mathcal{C}$ **do**
3:     $(M, X) \leftarrow \text{LOAD}(c)$; SetDeterministic()
4:     **for all** $(b_{\text{ref}}, b_{\text{tgt}}) \in B \times B$ **do**
5:         $Y_{\text{ref}} \leftarrow f(M, X; b_{\text{ref}})$; $Y_{\text{tgt}} \leftarrow f(M, X; b_{\text{tgt}})$
6:         $s_{\text{tensor}} \leftarrow \text{COMPARETENSOR}(Y_{\text{ref}}, Y_{\text{tgt}})$
7:         $s_{\text{task}} \leftarrow \text{COMPARETASK}(\cdot)$; LogJsonl($\cdot$)
8:     **end for**
9: **end for**

---

## IV. Experimental Setup

**Models.** ResNet18/50, MobileNetV3, ViT-B/16, Faster R-CNN, RetinaNet, YOLOv5n, UNet, DeepLabV3, FCN-ResNet50.
**Backends.** CPU (eager), GPU (eager), Compiled (`torch.compile`); optional ONNX Runtime / TensorRT.
**Tolerances.** $\text{atol} \in \{10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}\}$, $\text{rtol} = 10^{-5}$.
**Inputs.** Public-domain or synthetic; fixed preprocessing (resize/interp/normalize).
**Hardware.** We log GPU/driver/CUDA/cuDNN versions, CPU, RAM, and seeds/determinism flags to bound extraneous variability.

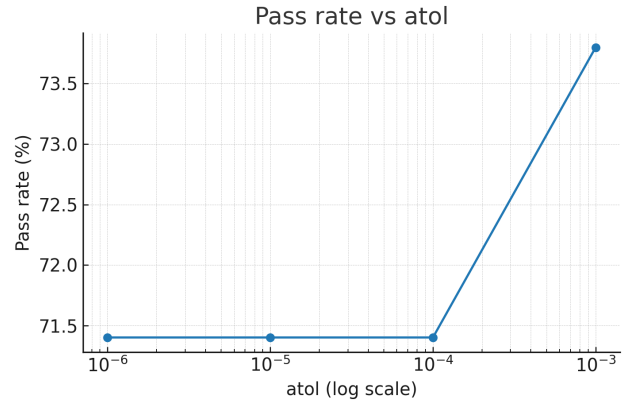## V. Results: Answers to RQ1–RQ3

### A. RQ1: Drift vs. Tolerances



Fig. 3: Pass rate across atol values.

TABLE I: Representative prior work versus this paper.

| Work | Primary Target | Granularity | Cross-backend drift | Deterministic adapters | Task metrics link |
|---|---|---|---|---|---|
| DeepXplore [1] | Testing | Input | ○ | ○ | ○ |
| DeepTest [2] | Testing | Input | ○ | ○ | ✓ |
| TensorFuzz [3] | Testing | Input/Interface | ○ | ○ | ○ |
| Mist [4] | Testing | Multi-API | ○ | ○ | ○ |
| TVM/XLA/Glow/Inductor [5] | Compiler | Graph/Kernels | △ | ○ | ○ |
| **Ours** | **Compatibility** | **Tensor/Task** | ✓ | ✓ | ✓ |

Legend: ✓ explicit; △ partial/indirect; ○ not a primary focus.

TABLE II: Threshold sensitivity: pass rate by atol.

| atol | Total | Passed | Pass % |
|---|---|---|---|
| $1e{-}6$ | 168 | 120 | 71.4 |
| $1e{-}5$ | 168 | 120 | 71.4 |
| $1e{-}4$ | 168 | 120 | 71.4 |
| $1e{-}3$ | 168 | 124 | 73.8 |

*Finding.* Pass rates improve monotonically as atol relaxes. Most failures concentrate at $10^{-6}$, indicating that fine-grained numerical perturbations are the principal driver; deployments should calibrate thresholds to task sensitivity.



Fig. 5: Failure taxonomy distribution.
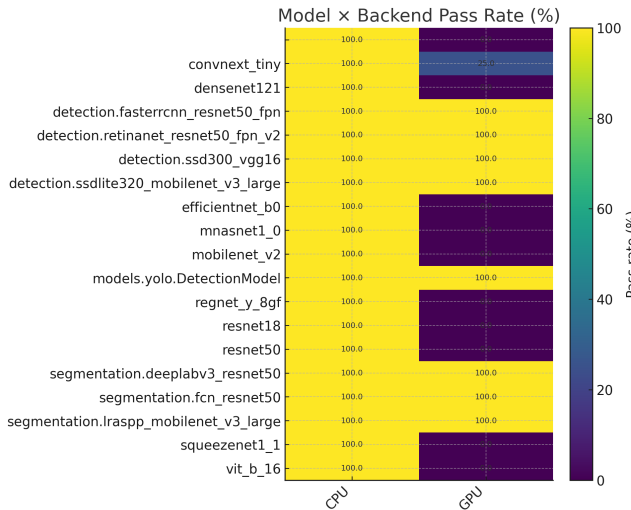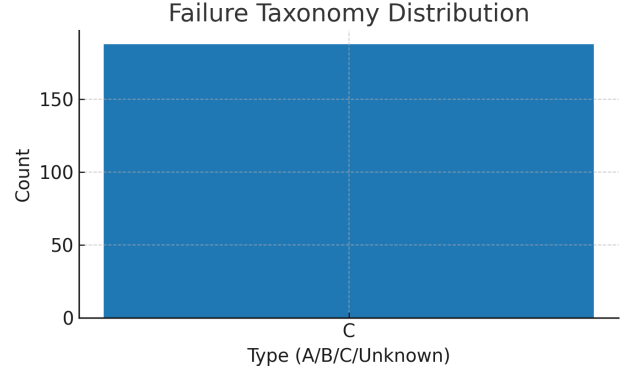
*Finding.* Detection models show lower agreement on compiled backends; the taxonomy indicates that ordering/tie-breaking in post-processing and partial operator support dominate failures. Segmentation tasks are comparatively stable, likely due to fewer order-sensitive operations.

### B. RQ2: Which Models/Tasks Diverge Most?

### C. RQ3: Where Does Divergence Emerge?

*Finding.* Selective activation probes suggest that early convolutional layers can seed drift for classification models, with discrepancies compounding in detection heads. A full-scale activation survey remains future work.

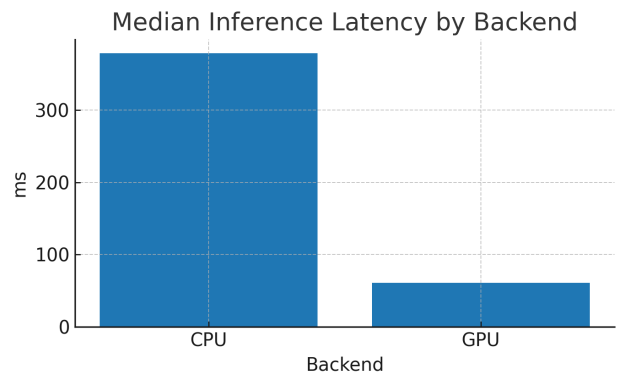### D. Latency and Trade-offs



Fig. 6: Median inference latency by backend.



Fig. 4: Pass-rate heatmap over Model × Target Backend.

*Finding.* Compiled backends reduce median latency for several architectures, illustrating an accuracy–performance trade-off

when compatibility gaps appear—underscoring the value of targeted stabilizers (deterministic adapters, fallbacks).

*E. Overall Summary*

TABLE III: Overall experiment summary.

| Metric | Value | Notes |
|---|---|---|
| Total checks | 672 | four atol settings |
| Passed | 484 | aggregate across models/backends |
| Pass rate | 72.0% | overall |
| Distinct models | 19 | classification/detection/segmentation |
| Target backends | 2 | GPU (eager), compiled |
| Distinct atol | 4 | $\{10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}\}$ |

## VI. CASE STUDY: DETECTION DRIFT FROM NONDETERMINISTIC NMS

**Symptom.** CPU vs. GPU/compiled runs exhibit inconsistent pre-NMS box ordering, leading to task-level discrepancies despite small tensor-level differences.

**Observation.** Agreement remains high up to the detection head; deviations emerge at the pre-NMS ordering stage, consistent with nondeterministic tie-breaking rather than upstream feature misalignment.

**Mitigation.** Deterministic sort over $(\text{score}, x_1, y_1)$ prior to NMS; alternatively, enforce FP32 for unstable kernels.

**Re-validation.** At $\text{atol} = 10^{-5}$, deterministic sorting restores task-level agreement without degrading latency benefits from compilation.



Before deterministic sort          After deterministic sort

Fig. 7: Qualitative detection comparison on the same image: left shows inconsistent NMS outcomes across backends; right shows alignment after enforcing deterministic pre-NMS sorting.

## VII. DISCUSSION

**Threats to Validity.** Hardware/driver autotuning, precision modes (AMP vs. FP32), preprocessing mismatches (resize/interp), and batch-size effects may introduce residual nondeterminism. While we log environment fingerprints, remaining variance cannot be fully excluded.

**Lessons.**

- **Tolerance calibration.** Very tight thresholds ($10^{-6}$) surface numerically small yet order-sensitive perturbations; thresholds should reflect end-task tolerance.
- **Deterministic adapters.** Sorting candidates before NMS removes order-induced divergence in detection with negligible overhead.
- **Operator fallbacks.** For problematic kernels, selective eager/FP32 fallbacks improve stability while preserving most performance gains.

**Future Work.** We plan a systematic activation-level survey across architectures and backends; broader model families including generative and multimodal; and additional runtimes (ONNX Runtime, TensorRT). We also aim to integrate the framework with reliable, efficient *foundation models*, aligning with emerging research priorities in reproducible, cross-platform deployment.

## VIII. REPRODUCIBILITY AND ARTIFACTS

We provide sanitized code, environment lockfiles, and scripts to regenerate JSONL logs and all tables/figures. An anonymized artifact is available for review and will be released publicly upon acceptance. The updated implementation is available at https://github.com/william-zehua-li/cross-backend-model-checker. All experiments can be reproduced with the provided configs and scripts.

## IX. CONCLUSION

We introduced a configuration-first framework for assessing cross-backend compatibility with a three-tier verifier that links tensor-level drift to task-level outcomes. Across 672 checks spanning four tolerance settings, 72.0% of runs pass; enforcing deterministic pre-NMS sorting restores detection-level agreement without forfeiting the latency benefits of compilation. *To our knowledge, this is the first unified framework that systematically quantifies and mitigates cross-backend drift under a common configuration-first protocol.* We believe this advances dependable deployment of deep learning in safety-critical domains—such as medical imaging and autonomous systems—where cross-backend consistency is essential for reproducibility and assurance.

## REFERENCES

[1] K. Pei, Y. Cao, J. Yang, and S. Jana, "DeepXplore: Automated Whitebox Testing of Deep Learning Systems," in *Proceedings of the 26th ACM Symposium on Operating Systems Principles (SOSP)*, 2017.

[2] Y. Tian, K. Pei, S. Jana, and B. Ray, "DeepTest: Automated Testing of Deep-Neural-Network-Driven Autonomous Cars," in *Proceedings of the 40th International Conference on Software Engineering (ICSE)*, 2018.

[3] N. Odena, C. Olsson, D. Andersen, and I. Goodfellow, "Tensor-Fuzz: Debugging Neural Networks with Coverage-Guided Fuzzing," *arXiv:1807.10875*, 2018.

[4] J. Zhang, X. Zhang, Y. Wei, et al. , "Mist: Automated Neural Network Model Testing via Multiple Interfaces," in *Proceedings of the 35th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, 2020.

[5] T. Chen, T. Moreau, Z. Jiang, et al. , "TVM: An Automated End-to-End Optimizing Compiler for Deep Learning," in *Proceedings of the 13th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, 2018. (See also: OpenXLA/XLA; Glow; PyTorch 2.0 Inductor.)

[6] J. Pineau, P. Vincent-Lamarre, K. Sinha, et al. , "Improving Reproducibility in Machine Learning Research: A Report from the NeurIPS 2019 Reproducibility Program," *Journal of Machine Learning Research*, 22(164):1–20, 2021.

## APPENDIX

```
1  # Library example (classification)
2  from: library
3  model: resnet18
4  inputs: [assets/cat.jpg]
5  means: [0.485, 0.456, 0.406]
6  stds: [0.229, 0.224, 0.225]
7  options: { compile: true, resize_multiple:
       32 }
8  verification: { tol: { atol: 1e-5, rtol: 1e
       -5 } }
```

```
1  # Repo example (local clone + class path)
2  from: repo
3  repo: { path: ../third_party/yolo_clone,
       class: models.yolo.YoloNet }
4  params: { img_size: 640 }
5  inputs: [data/dog.jpg]
```

```python
1  #!/usr/bin/env python3
2  import argparse, glob, os
3  from typing import Any, List, Tuple, Union
4  import torch
5  from PIL import Image
6  from sanitized_utils import (
7      set_global_seed, load_yaml, resolve_path,
           load_image_tensor,
8      adjust_to_multiple, tensors_allclose,
           to_cpu_like,
9  )
10 from sanitized_loaders import LibraryLoader,
        RepoLoader
11 TensorOrList = Union[torch.Tensor, List[
       torch.Tensor]]
12
13 def _normalize_output(y: Any) -> List[torch.
       Tensor]:
14     if isinstance(y, torch.Tensor): return [y
           ]
15     if isinstance(y, (list, tuple)): return [
           t for t in y if isinstance(t, torch.
           Tensor)]
16     if isinstance(y, dict): return [v for v
           in y.values() if isinstance(v, torch.
           Tensor)]
17     return []
18
19 def run_once(cfg_path: str, device: str,
       use_compile: bool) -> Tuple[TensorOrList
       , TensorOrList]:
20     cfg = load_yaml(cfg_path)
21     source = cfg.get("from", "library")
22     means = cfg.get("means", [0.485, 0.456,
           0.406])
23     stds = cfg.get("stds", [0.229, 0.224,
           0.225])
24     inputs = cfg.get("inputs", [])
25     options = cfg.get("options", {})
26     resize_multiple = options.get("
           resize_multiple", 32)
27     if source == "library":
28         model_name = cfg.get("model")
29         if not model_name: raise ValueError("
               For 'library' source you must
               specify 'model'.")
30         loader = LibraryLoader(model_name=
               model_name, weights=cfg.get("
               weights"))
31         model = loader.build()
32     elif source == "repo":
33         repo = cfg.get("repo", {})
34         repo_path = resolve_path(cfg_path,
               repo.get("path", ""))
35         class_path = repo.get("class", "")
36         params = cfg.get("params", {})
37         loader = RepoLoader(repo_path=
               repo_path, class_path=class_path,
               params=params)
38         model = loader.build()
39     else:
40         raise ValueError(f"Unknown source: {
               source}")
41     model.eval()
42     ref_device = torch.device("cpu")
43     model_ref = model.to(ref_device)
44     tgt_device = torch.device(device)
45     model_tgt = model.to(tgt_device)
46     if use_compile: model_tgt = torch.compile
           (model_tgt)
47     ref_outputs: List[torch.Tensor] = []
48     tgt_outputs: List[torch.Tensor] = []
49     for rel in inputs:
50         img_path = resolve_path(cfg_path, rel)
51         img = Image.open(img_path).convert("
               RGB")
52         x = load_image_tensor(img, means, stds
               ).unsqueeze(0)
53         x_ref = x.to(ref_device); x_tgt = x.to
               (tgt_device)
54         if resize_multiple:
55             x_ref = adjust_to_multiple(x_ref,
                   resize_multiple)
56             x_tgt = adjust_to_multiple(x_tgt,
                   resize_multiple)
57         with torch.no_grad():
58             y_ref = model_ref(x_ref); y_tgt =
                   model_tgt(x_tgt)
59         ref_outputs.extend([to_cpu_like(t) for
                t in _normalize_output(y_ref)])
60         tgt_outputs.extend([to_cpu_like(t) for
                t in _normalize_output(y_tgt)])
61     return (ref_outputs if len(ref_outputs)
           != 1 else ref_outputs[0],
62         tgt_outputs if len(tgt_outputs) !=
               1 else tgt_outputs[0])
63
64 def main():
65     parser = argparse.ArgumentParser(
           description="Sanitized compatibility
           test runner")
66     parser.add_argument("-d", "--device",
           required=True, help="Target device, e.
           g. cpu or cuda")
67     parser.add_argument("-c", "--configs",
           default="configs/*.yaml", help="Glob
           for YAML configs")
68     parser.add_argument("--compile", action="
```

```python
                   store_true", help="Use torch.compile
                       for target run")
69     parser.add_argument("--seed", type=int,
           default=5, help="Global RNG seed")
70     args = parser.parse_args()
71     set_global_seed(args.seed)
72     cfg_files = sorted(glob.glob(args.configs
           ))
73     if not cfg_files:
74         print(f"No configs matched: {args.
               configs}"); return
75     total, passed, failed = 0, 0, 0
76     for cfg in cfg_files:
77         total += 1
78         try:
79             ref, tgt = run_once(cfg, args.
                   device, args.compile)
80             conf = load_yaml(cfg)
81             tol = (((conf.get("verification")
                   or {}).get("tol")) or {})
82             atol = float(tol.get("atol", 1e-5))
                   ; rtol = float(tol.get("rtol",
                   1e-5))
83             ok = tensors_allclose(ref, tgt,
                   atol=atol, rtol=rtol)
84             status = "PASS" if ok else "FAIL"
85             if ok: passed += 1
86             else: failed += 1
87             print(f"[{status}] {os.path.
                   basename(cfg)} (atol={atol},
                   rtol={rtol})")
88         except Exception as e:
89             failed += 1
90             print(f"[ERROR] {os.path.basename(
                   cfg)} -> {e}")
91     print("\n=== Summary ==="); print(f"Total
           : {total} Passed: {passed} Failed: {
           failed}")
92
93 if __name__ == "__main__": main()
```

```python
16         if isinstance(state, dict) and "
               state_dict" in state: state =
               state["state_dict"]
17         model.load_state_dict(state, strict
               =False)
18     return model
19
20 class RepoLoader:
21     def __init__(self, repo_path: str,
           class_path: str, params: Optional[
           Dict[str, Any]] = None):
22         if not repo_path or not os.path.isdir(
               repo_path): raise
               FileNotFoundError(f"repo_path not
               found: {repo_path}")
23         if "." not in class_path: raise
               ValueError("class_path must be
               dotted, e.g. 'pkg.subpkg.Class'")
24         self.repo_path = os.path.abspath(
               repo_path); self.class_path =
               class_path; self.params = params
               or {}
25     def build(self) -> nn.Module:
26         import sys
27         sys.path.insert(0, self.repo_path)
28         try:
29             module_path, cls_name = self.
                   class_path.rsplit(".", 1)
30             module = importlib.import_module(
                   module_path); cls = getattr(
                   module, cls_name)
31             model = cls(**self.params); return
                   model
32         finally:
33             if self.repo_path in sys.path: sys.
                   path.remove(self.repo_path)
```

```python
1  # Sanitized loaders: only public sources; no
       proprietary modules.
2  import importlib, os
3  from typing import Any, Dict, Optional
4  import torch, torch.nn as nn
5
6  class LibraryLoader:
7      def __init__(self, model_name: str,
           weights: Optional[str] = None, params
           : Optional[Dict[str, Any]] = None):
8          self.model_name = model_name; self.
               weights = weights; self.params =
               params or {}
9      def build(self) -> nn.Module:
10         from torchvision import models
11         if not hasattr(models, self.model_name
               ):
12             raise ValueError(f"Unknown library
                   model: {self.model_name}")
13         ctor = getattr(models, self.model_name
               ); model = ctor(**self.params)
14         if self.weights and os.path.exists(
               self.weights):
15             state = torch.load(self.weights,
                   map_location="cpu")
```

```python
1  # Utilities for the sanitized runner.
2  import os, random
3  from typing import List, Sequence, Union
4  import numpy as np, torch
5  from PIL import Image
6  import yaml
7  from torchvision import transforms
8
9  def set_global_seed(seed: int) -> None:
10     os.environ["PYTHONHASHSEED"] = str(seed);
           random.seed(seed); np.random.seed(
           seed)
11     torch.manual_seed(seed); torch.cuda.
           manual_seed_all(seed)
12     torch.backends.cudnn.deterministic = True
           ; torch.backends.cudnn.benchmark =
           False
13
14 def load_yaml(path: str) -> dict:
15     with open(path, "r", encoding="utf-8") as
           f: return yaml.safe_load(f) or {}
16
17 def resolve_path(base_cfg: str, rel: str) ->
       str:
18     if os.path.isabs(rel): return rel
19     base_dir = os.path.dirname(os.path.
           abspath(base_cfg)); return os.path.
           normpath(os.path.join(base_dir, rel))
20
```

```python
21  def load_image_tensor(img: Image.Image,
        means: Sequence[float], stds: Sequence[
        float]) -> torch.Tensor:
22      pre = transforms.Compose([transforms.
            Resize(256), transforms.CenterCrop
            (224), transforms.ToTensor(),
23                          transforms.Normalize(
                                mean=list(means),
                                 std=list(stds))
                                ,])
24      return pre(img)

25
26  def adjust_to_multiple(x: torch.Tensor, m:
        int) -> torch.Tensor:
27      if x.dim() != 4: return x
28      _, _, h, w = x.shape; nh = max(m, (h // m
            ) * m); nw = max(m, (w // m) * m)
29      if nh == h and nw == w: return x
30      return torch.nn.functional.interpolate(x,
            size=(nh, nw), mode="bilinear",
            align_corners=False)

31
32  def to_cpu_like(t: torch.Tensor) -> torch.
        Tensor: return t.detach().to("cpu")

33
34  def _allclose(a: torch.Tensor, b: torch.
        Tensor, atol: float, rtol: float) ->
        bool:
35      try: torch.testing.assert_close(a, b,
            atol=atol, rtol=rtol); return True
36      except AssertionError: return False

37
38  def tensors_allclose(a: Union[torch.Tensor,
        List[torch.Tensor]], b: Union[torch.
        Tensor, List[torch.Tensor]], atol: float
         = 1e-5, rtol: float = 1e-5) -> bool:
39      if isinstance(a, torch.Tensor) and
            isinstance(b, torch.Tensor): return
             _allclose(a, b, atol, rtol)
40      if isinstance(a, list) and isinstance(b,
             list):
41          if len(a) != len(b): return False
42          for ta, tb in zip(a, b):
43              if not _allclose(ta, tb, atol, rtol
                    ): return False
44          return True
45      return False
```

```python
        backend_pair: str,
13          atol: float, rtol: float, status:
                str, stats: Dict[str, Any],
                out_path: str):
14      record = {
15          "config": config_path, "model": model,
                "backend_pair": backend_pair,
16          "atol": atol, "rtol": rtol, "status":
                status, **(stats or {})
17      }
18      log_jsonl(record, out_path)
```

```python
1   # Minimal JSONL logger to unify outputs for
        paper tables/figures.
2   import json, time, os, sys
3   from typing import Any, Dict

4
5   def log_jsonl(record: Dict[str, Any], path:
        str) -> None:
6       rec = dict(record)
7       rec["timestamp"] = time.strftime("%Y-%m-%
            dT%H:%M:%S")
8       os.makedirs(os.path.dirname(path),
            exist_ok=True)
9       with open(path, "a", encoding="utf-8") as
             f:
10          f.write(json.dumps(rec, ensure_ascii=
                False) + "\n")

11
12  def log_run(config_path: str, model: str,
```