

What Theory Does Deep Learning Need: A Dynamics and Measurement Perspective

Responding to Yuandong Tian's question: what would a consensus mathematical theory of deep learning look like in 100 years?

TL;DR

If deep learning ever gets a widely accepted "mathematical theory" that explains today's training "magic", it will probably look less like a closed axiomatic system about `theta` and `-grad L` and more like a **dynamical + measurement theory**:

- **Dynamics:** training as a stochastic process with regime changes (phase-like behavior), not a single smooth curve.
- **Constraints:** architecture, data, compute, and governance boundaries as first-class objects that shape what is reachable.
- **Time:** "tempo" (update and feedback time scales) as a control variable, not a background parameter.
- **Measurement:** explicit estimators and audit gates, so claims are falsifiable rather than post-hoc narratives.

Without the measurement layer, every "deep learning theory" is just a retroactive story.

This post is a research-and-engineering agenda: a way to turn "training intuition" into auditable objects you can measure, test, and (if needed) falsify.

1) The object is not just parameters: it's a system state

Many "training mysteries" are symptoms of analyzing a multi-variable system with a one-variable lens. A future theory will not treat the model as just a parameter vector. It will treat training as an evolving **system state** that includes optimization state, data process, evaluation gates, and other boundaries:

$$x(t) = (\text{theta}(t), s(t), D(t), \text{Pi}(t), B(t), C(t))$$

One possible reading:

- `theta(t)`: parameters
- `s(t)`: optimizer and training-internal state (moments, noise scale proxies, etc.)
- `D(t)`: data process (mixtures, curricula, filters, licensing boundaries)
- `Pi(t)`: representation geometry proxies (effective rank, spectra, sparsity, etc.)
- `B(t)`: boundary conditions (training budget, deployment gate cadence, human oversight cycle)
- `C(t)`: constraint structure (how the reachable space is being shaped over time)

Once you write the object this way, "magic" becomes hypotheses about missing variables.

2) Training is not a single gradient: it's a resultant force plus structured noise

At a high level, training behaves more like a stochastic dynamical system than a pure gradient flow:

$$dx = f(x; B, C) dt + \Sigma(x; B, C) dW_t$$

- f aggregates many “forces”: task gradients, regularization, curriculum effects, alignment losses, architectural biases, tool-use constraints, and so on.
- Σ captures structured noise: SGD sampling noise, distributed training artifacts, data noise, and changes in sampling strategy.

Many “hyperparameter tricks” are interventions on the shape of f or the spectrum of Σ .

3) Constraints are the real “magic-killer”

Deep learning discussions often treat constraints as engineering details. This is a mistake. A mature theory will elevate constraints to first-class theoretical objects because they define **reachability**—what endpoints are even possible given your setup.

Constraints come from:

- Architecture constraints (attention structure, inductive biases, available tools)
- Data constraints (coverage, long-tail, noise, licensing, filtering)
- Compute constraints (budget, parallelism, communication noise)
- Governance constraints (evaluation gates, approval cadence, rollback capacity)

A core design lesson (validated in simple dynamical testbeds) is:

Boundary conditions are part of the constraint structure. They define which endpoints are even reachable.

In practice, “the same training code” under different boundaries (data mixture constraints, evaluation gates, rollout boundaries) can converge to qualitatively different regimes.

4) Time is not a background parameter: regimes and tempo matter

The phenomena practitioners call “magic” often look like regime switching:

- Early phase: rapid fitting (optimization-dominated)
- Middle phase: representation reorganization (structure changes)
- Late phase: stabilization / compression / attractors (slow drift, high rigidity)

A consensus theory will not try to explain every point. It will aim to:

- Detect regime changes with observable signatures
- Predict how “phase boundaries” move as constraints change

- Explain why different systems share similar behaviors under similar macro-conditions (universality classes)
-

5) The missing layer: measurement (estimators) and audit gates

Here is where many “big perspectives” fail: they remain descriptive because they do not specify how to measure what they invoke.

If we want a theory that is falsifiable rather than narrative, we need an explicit measurement layer: statements bound to a declared estimator tuple:

$$E = (S_t, B, \{F_{\hat{t}}, C_{\hat{t}}, I_{\hat{t}}\}, W)$$

- S_t : what you treat as the “state” (parameter trajectory, activations, logs, etc.)
- B : boundary conditions (data distribution, rollout boundary, gating cadence)
- $\{F_{\hat{t}}, C_{\hat{t}}, I_{\hat{t}}\}$: operational estimators (force proxies, constraint proxies, information proxies)
- W : measurement window / hyperparameters

And we need at least one minimal audit concept:

If you claim “constraint is increasing”, you should show that multiple reasonable $C_{\hat{t}}$ estimators agree on the task-relevant structure (rank order, thresholds, or event structure). Otherwise you may be measuring an artifact.

This kind of coherence gate is not philosophy. It’s what makes cross-domain statements testable.

6) Tier-2 validation can’t be “relabeling”: an A/B/C novelty filter

A good warning for any meta-framework is: you can relabel old ideas with new words and call it progress.

To avoid that, here is a hard gate for any Tier-2 (real-world) case study:

- **A**: How does the field describe the risk today (without this framework)?
- **B**: How does the framework translate that description into a state / force / constraint / tempo story?
- **C**: What new, quantitative, falsifiable prediction follows that is not already contained in (A)?

If you can't answer (C), the case study is a demo, not a scientific validation. Most “framework applications” fail this test.

7) Three concrete, falsifiable Tier-2 predictions (examples)

These are “shape of curve” predictions, not slogans. They can be wrong.

T2-P1: A tempo threshold where measurement coherence collapses

Hypothesis: as update tempo increases, you cross a critical point where different “reasonable” estimators stop agreeing, signaling that the system has entered a regime where governance and measurement cannot keep up.

One operational sketch:

- There exists a τ_c such that coherence $\rho(\hat{C}_1, \hat{C}_2)$ drops sharply once $\tau_u >= \tau_c$.
- τ_u : update interval (training or deployment changes per unit time)
- $\rho(\hat{C}_1, \hat{C}_2)$: estimator coherence (a Spearman-like proxy)

Prediction: beyond τ_c you should see coherence degrade sharply, before obvious failures occur.

T2-P2: IO event rate changes shape under sustained tempo mismatch

Hypothesis: under sustained mismatch between evaluation cycle and update cycle, irreversible operations (IOs) accumulate nonlinearly.

- $\text{IO_rate}(t)$ is approximately linear when $\tau_g / \tau_u \leq r$, and roughly exponential when $\tau_g / \tau_u > r$.
- τ_g : governance/evaluation cycle time
- r : mismatch threshold (domain-specific)

Prediction: there exists a regime where “nothing obviously breaks” yet irreversibility indicators accelerate.

T2-P3: Boundary changes decide endpoints (boundary = constraint)

Hypothesis: changing boundary conditions (data boundaries, rollout boundaries, gate cadences) shifts the reachable regime structure and thus the likely endpoints, even if you keep the “core training algorithm” constant.

This is testable by controlled comparisons: same training recipe, same compute, different boundary constraints.

8) A two-week pilot for practitioners (the lowest-friction bridge)

If you want traction with serious ML engineering leaders, give them a low-cost pilot they can run without adopting a worldview.

8.1 Three dashboard metrics

1. **Validation Lag (VL)**: time from change-effective to evaluation closure
2. **Rollback Drill Pass Rate (RDPR)**: fraction of rollback rehearsals that succeed within a defined RTO/RPO
3. **Gate Bypass Rate (GBR)**: rate of bypassing required gates for high-impact changes

8.2 A minimal IO register (classification, not blame)

An IO is a change that permanently shrinks feasible future correction pathways under bounded cost/time.

Minimal categories:

- Data IO (non-reproducible sources, irreversible filtering, major mixture shifts)
- Evaluation IO (shortening/removing gates, allowing fast paths)
- Alignment IO (policy/execution changes that reduce auditability or override capacity)
- Deployment IO (expanding exposure scope, removing staging boundaries)
- Supply-chain IO (non-auditable components in critical paths)
- Optionality IO (removing redundancy, single-point dependency lock-in)

Pilot deliverable: label the last 30 to 90 days of high-impact changes and backtest the three metrics against incidents, rollbacks, and near misses.

9) Where this connects to FIT (and why I think it's useful)

The lens behind the structure above is the FIT framework (Force–Information–Time) plus an explicit estimator selection layer (EST).

You don't need to buy the "universal framework" claim to use the engineering outputs:

- a system-state view rather than parameter-only
- explicit boundary/constraint modeling
- estimator tuples and coherence gates
- Tier-2 novelty filters that force new predictions
- an IO register and minimal tempo governance pilot

If any of these fail under honest testing, that failure is informative.

Links (anchor)

- FIT repository (GitHub): <https://github.com/qienhuang/F-I-T>
 - FIT v2.4 specification (EST edition): <https://github.com/qienhuang/F-I-T/blob/main/docs/v2.4.md>
 - EST quickstart notebook: https://github.com/qienhuang/F-I-T/blob/main/experiments/est_quickstart/fit-quickstart.ipynb
 - Tier-2 predictions register (ChinaXiv package): https://github.com/qienhuang/F-I-T/blob/main/papers/chinaxiv/chinaxiv_tier2_predictions_register.zh_cn.md
 - Zenodo archive (v2.4 release): <https://doi.org/10.5281/zenodo.18082325>
 - Gitee mirror: <https://gitee.com/qienhuang/f-i-t>
-

Disclosure

Drafting and editing were assisted by language models. The author takes responsibility for all claims and errors.