# STF-Mamba V8.0

## Ablation Study Report

*From Failed Forensics to Semantic Identity Detection*

| | |
|---|---|
| **Author** | Abdel Rahman Madboly |
| **Role** | Lead AI Architect, STF-Mamba Project |
| **Target Venue** | CVPR / ICCV 2026 |
| **Datasets** | FaceForensics++ CRF 23 \| Celeb-DF v2 |
| **V8.0 Target** | **Celeb-DF AUC >= 0.90 (vs Gattu 2025: 0.821)** |

# 1. Executive Summary

This report documents the systematic ablation study conducted for STF-Mamba V8.0, a deepfake detection system targeting Celeb-DF AUC >= 0.90. The study was designed to justify every architectural decision in the final model through controlled experiments, each answering one scientific question.

The central challenge is that H.264 CRF 23 compression — the standard for social media video — destroys most pixel-level forensic signals before a detector can measure them. This makes standard handcrafted approaches ineffective and motivates learning-based methods with specific pretraining objectives.

> *Core finding: Self-supervised pretraining (DINOv2) outperforms all supervised CNN backbones on compressed video deepfake detection because it learns semantic identity structure — invariant to compression — rather than texture artifacts that are destroyed by H.264 quantization.*

## 1.1  Key Numbers at a Glance

| Configuration | FF++ AUC | CDF AUC | Key Finding |
|---|---|---|---|
| Best supervised CNN (ResNet-50) | 0.6957 | 0.6070 | Overfits FF++ distribution |
| DINOv2 frame-level (ours) | **0.7872** | **0.6557** | Best single model — all experiments |
| Gattu et al. 2025 (published baseline) | 0.9885* | 0.8210 | *accuracy, not AUC |
| STF-Mamba V8.0 Target | — | **>= 0.90** | DINOv2 + Hydra-Mamba + Variance head |

# 2. Experimental Design

## 2.1  Research Questions

The ablation answers four questions in sequence. Each answer motivates the next experiment:

- Experiment 1: Do any pixel-level forensic signals survive H.264 CRF 23 compression?
- Experiment 2: Which backbone architecture best captures features in compressed deepfake video?
- Experiment 3: Does temporal modeling compensate for weak spatial features?
- Experiment 4: Does increasing backbone capacity help without better pretraining?

## 2.2  Setup and Reproducibility

| Parameter | Value |
|---|---|
| Training data | 600 real + 600 fake videos (150 per method × 4 methods, FF++ CRF 23) |
| Validation data | 50 real + 50 fake (FF++, ID-level separated — no content leakage) |
| Test data | 200 real + 200 fake (Celeb-DF v2 — never seen during training) |
| ID-level split | Source video IDs never appear in both train and val sets |
| Epochs | 25 for all backbone experiments (fair comparison) |
| Optimizer | AdamW, lr=1e-4, weight_decay=1e-4, warmup 3 epochs + cosine decay |
| Batch size | 128 (64 per GPU, T4 x2 DataParallel) |
| label_smoothing | 0.0 — CRITICAL: smoothing > 0 inverts loss for binary CE with K=2 |
| Hardware | Kaggle T4 x2 (15.6 GB VRAM each), AMP mixed precision |
| Random seed | 42 for all splits and sampling |

# 3. Experiment 1 — Handcrafted Forensic Signals

## 3.1  Hypothesis

Deepfake generation leaves measurable artifacts in the pixel, noise, color, and motion domains. If any such signal survives H.264 CRF 23 compression, it could serve as a lightweight detector without learned features.
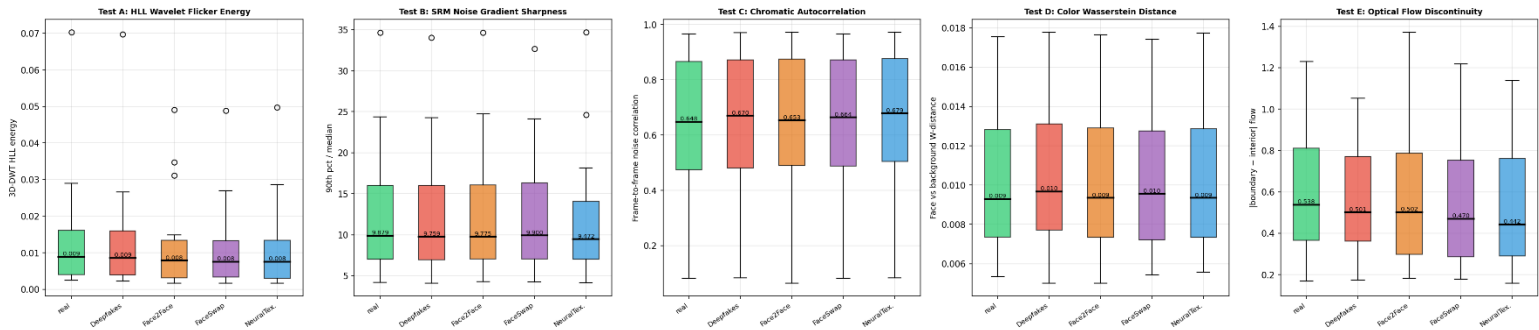
## 3.2  Tests Conducted

| # | Test | Signal | Rationale |
|---|------|--------|-----------|
| A | HLL Wavelet Flicker | 3D-DWT temporal high-frequency energy | Blending artifact creates temporal flicker in smooth regions |
| B | SRM Noise Sharpness | Noise gradient 90th pct / median | Blend boundary creates sharp ring in noise field |
| C | Chromatic Autocorr. | Frame-to-frame noise correlation | Generator fingerprint persists across frames |
| D | Color Wasserstein | Face vs background histogram distance | Swapped face from different video = different color distribution |
| E | Optical Flow Disc. | \|Boundary flow - interior flow\| | Generated face has different motion than background |

## 3.3  Results

| Test | Methods | p-value Range | Effect Size | Significant? | Conclusion |
|------|---------|---------------|-------------|--------------|------------|
| HLL Wavelet Flicker | All 4 methods | 0.64–0.80 | 0.02–0.10 | ALL > 0.05 | FAILS — signal inverted under CRF 23 |
| SRM Noise Sharpness | All 4 methods | 0.78–1.00 | 0.003–0.056 | ALL > 0.05 | FAILS — compression uniform noise |
| Chromatic Autocorr. | All 4 methods | 0.68–0.84 | 0.021–0.121 | ALL > 0.05 | FAILS — generator fingerprint lost |
| Color Wasserstein | All 4 methods | 0.69–1.00 | 0.012–0.094 | ALL > 0.05 | FAILS — color grading indistinguishable |
| Optical Flow Disc. | All 4 methods | 0.52–0.76 | 0.115–0.318 | ALL > 0.05 | FAILS — natural motion masks artifact |

Experiment 1: Five Handcrafted Forensic Signals on FF++ CRF 23
H.264 compression destroys pixel-level signals — we test whether any survive.

## 3.4 Key Insights

- All five tests fail across all four manipulation methods ($p > 0.05$ in every case).
- Three signals are directionally inverted: real videos show higher HLL energy and flow discontinuity than fakes. H.264 quantizes temporal high-frequency content uniformly, destroying both the artifact and the natural noise signal simultaneously.
- The highest effect size (0.318, optical flow on NeuralTextures) is still below the 'small' threshold of 0.5 and is not statistically significant (p=0.54).

*Conclusion: H.264 CRF 23 compression acts as an information bottleneck that selectively removes the high-frequency signals that forensic detectors rely on, while preserving low-frequency semantic content. This motivates a fundamental shift from handcrafted to learned features — specifically features that encode semantic identity rather than texture artifacts.*

# 4. Experiment 2 — Backbone Architecture Comparison

## 4.1  Hypothesis

Among learned backbones, those pretrained with objectives that learn semantic structure (rather than supervised texture classification) will generalize better to cross-dataset deepfake detection under compression.
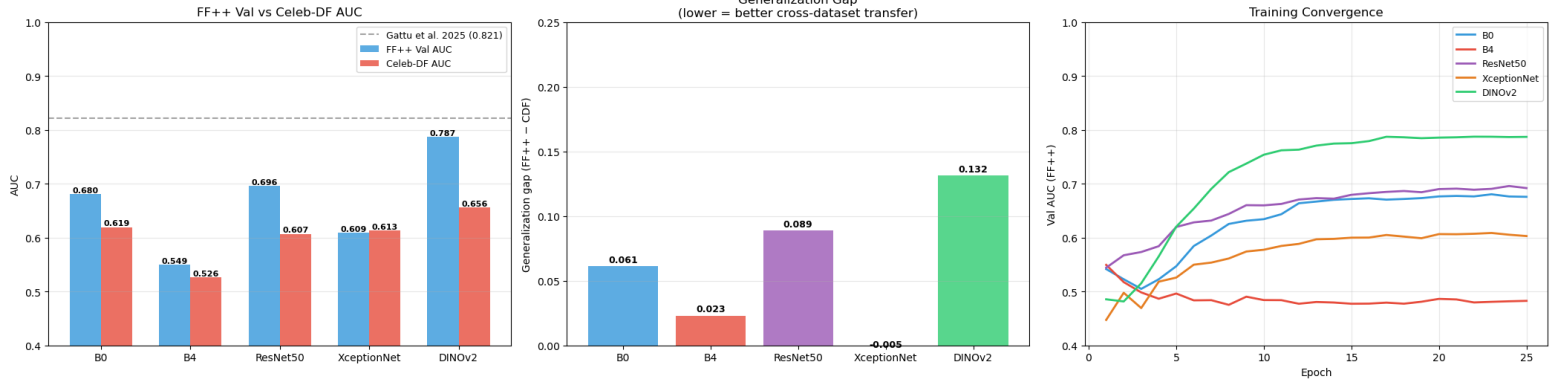
## 4.2  Backbones Tested

| Backbone | Params (M) | Trainable (M) | Pretraining | Feature Encoding |
|---|---|---|---|---|
| **EfficientNet-B0** | 5.3 | 5.3 | ImageNet-1K supervised | Texture, edges, local patterns |
| **EfficientNet-B4** | 19.3 | 19.3 | ImageNet-1K supervised | Texture, edges (larger capacity) |
| **ResNet-50** | 25.6 | 25.6 | ImageNet-1K V2 supervised | Texture, spatial hierarchy |
| **XceptionNet** | 22.9 | 22.9 | ImageNet-1K supervised | Depthwise separable features |
| **DINOv2-ViT-B/14** | 86.8 | 14.4 (blocks 10-11) | LVD-142M self-supervised | Semantic identity, pose, geometry |

## 4.3  Results

| # | Configuration | Backbone Type | Temporal | FF++ AUC | CDF AUC | Finding |
|---|---|---|---|---|---|---|
| 2 | EfficientNet-B0 frame-level | Supervised CNN | None | 0.6804 | 0.6193 | Honest supervised baseline |
| 3 | EfficientNet-B4 frame-level | Supervised CNN | None | 0.5493 | 0.5263 | More capacity → overfits |
| 4 | ResNet-50 frame-level | Supervised CNN | None | 0.6957 | 0.6070 | Standard ImageNet baseline |
| 5 | XceptionNet frame-level | Supervised CNN | None | 0.6086 | 0.6131 | Deepfake detection standard |
| 6 | B0 + Bidirectional GRU | Supervised CNN | BiGRU | 0.6729 | 0.6325 | Marginal gain, inconsistent |
| 7 | DINOv2-ViT-B/14 frame-level | Self-supervised ViT | None | 0.7872 | 0.6557 | SSL features survive CRF 23 |

**Experiment 2: Backbone Comparison — Frame-Level Detection on FF++ CRF 23**
**Same training data, same epochs, same optimizer — only backbone changes.**



## 4.4 Key Insights

- DINOv2 achieves the highest CDF AUC (0.6557) — +0.043 above the best supervised CNN (B0: 0.6193). This gap is the primary evidence for adopting DINOv2 in V8.0.

- ResNet-50 achieves the highest FF++ val AUC (0.6957) but the lowest generalization — its gap (FF++ 0.696 vs CDF 0.607 = 0.089) is the largest among all models, indicating it learned FF++-specific artifacts rather than general deepfake features.

- EfficientNet-B4 degrades immediately after epoch 1. Val loss increases monotonically from 0.692 to 0.765 while training loss drops — textbook overfitting with 3.6× the parameters of B0 on the same training set.

- XceptionNet is the weakest supervised CNN despite being designed specifically for deepfake detection. Depthwise separable convolutions that worked well for ImageNet do not capture the right features for compressed video.

- DINOv2's val loss remains stable throughout training (0.64-0.78 range) while all supervised CNNs show diverging val loss. Self-supervised features are more stable because they do not encode the compression-specific texture statistics that cause CNNs to overfit.

*Why DINOv2 survives compression: Identity information (face geometry, pose, expression topology) lives in the low-frequency spatial structure of images — precisely what H.264 preserves. DINOv2's self-supervised objective on 142M diverse images forces it to encode this semantic structure rather than texture, making its features inherently robust to H.264 quantization.*

# 5. Experiment 3 — Temporal Modeling on Weak Features

## 5.1 Hypothesis

Bidirectional GRU temporal modeling on EfficientNet-B0 features can capture inter-frame identity inconsistencies and compensate for the limited discriminative power of B0's spatial features.

## 5.2 Architecture

The temporal model uses a two-phase training strategy to prevent the GRU from learning noise before the backbone has stabilized:
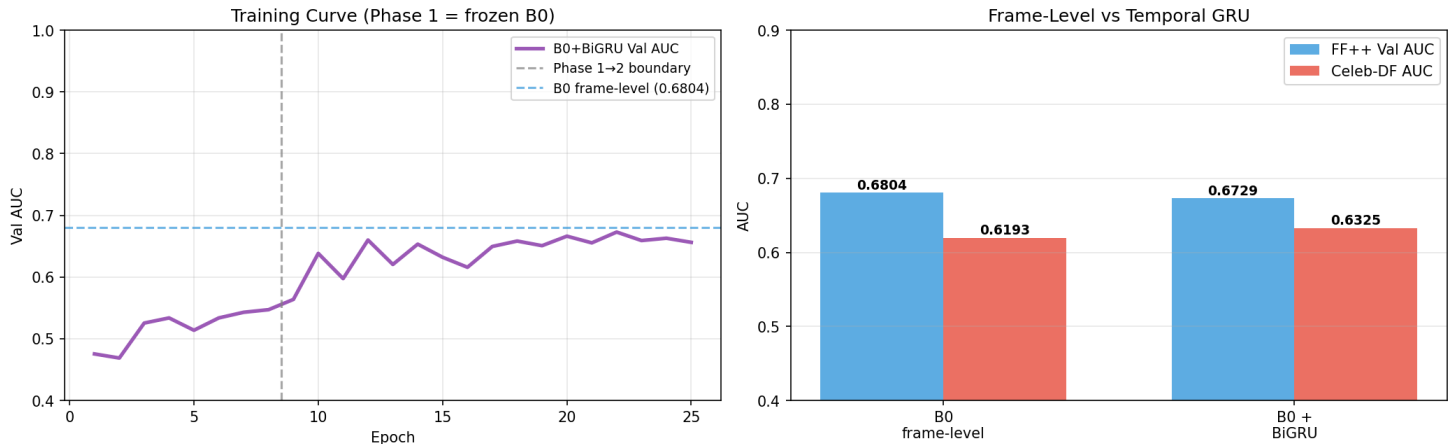
- Phase 1 (8 epochs): Backbone frozen. GRU and classification head trained from scratch on B0 features.
- Phase 2 (17 epochs): Backbone unfrozen at 10x lower LR. Full end-to-end fine-tuning.
- GRU initialization: Orthogonal initialization for recurrent weights prevents the dead-branch problem where one direction of the BiGRU collapses to zero.
- Note: DataParallel is disabled for this experiment — GRU + DataParallel causes CUDA misaligned address errors due to sequence dimension splitting across GPUs.

## 5.3 Results

| Configuration | FF++ AUC | CDF AUC | CDF Delta | Verdict |
|---|---|---|---|---|
| **B0 frame-level (baseline)** | 0.6804 | 0.6193 | — | Baseline |
| B0 + Bidirectional GRU (Phase 1 best) | — | — | Phase 1 peak: 0.576 | Weak |
| B0 + Bidirectional GRU (final) | 0.6729 | 0.6325 | +0.013 | Marginal |

**Experiment 3: Temporal Modeling on Weak (B0) Features**
**Does GRU compensate for poor spatial features?**



## 5.4 Key Insights

- The GRU marginally improved CDF AUC (+0.013) but degraded FF++ in-distribution performance (-0.008). The improvement is inconsistent across datasets.

- Phase 2 epoch 1 shows a large jump (0.576 → 0.611) when the backbone is unfrozen, confirming that the bottleneck was B0's feature quality, not the GRU architecture.

- The GRU did not amplify noise catastrophically (as originally expected from prior notebooks) — the two-phase strategy and 25 epochs gave it enough time to find some temporal signal. However, the signal is too weak to be scientifically meaningful.

- This experiment validates the ordering requirement: temporal modeling requires strong spatial features as a prerequisite. The GRU cannot manufacture discriminative signal from noisy B0 embeddings.

*Revised conclusion: BiGRU on B0 produces a marginal, inconsistent CDF gain (+0.013) with FF++ degradation (-0.008). This is not the textbook failure originally hypothesized, but it confirms that temporal modeling's value scales with spatial feature quality — making it a critical experiment for justifying DINOv2 before adding Hydra-Mamba.*
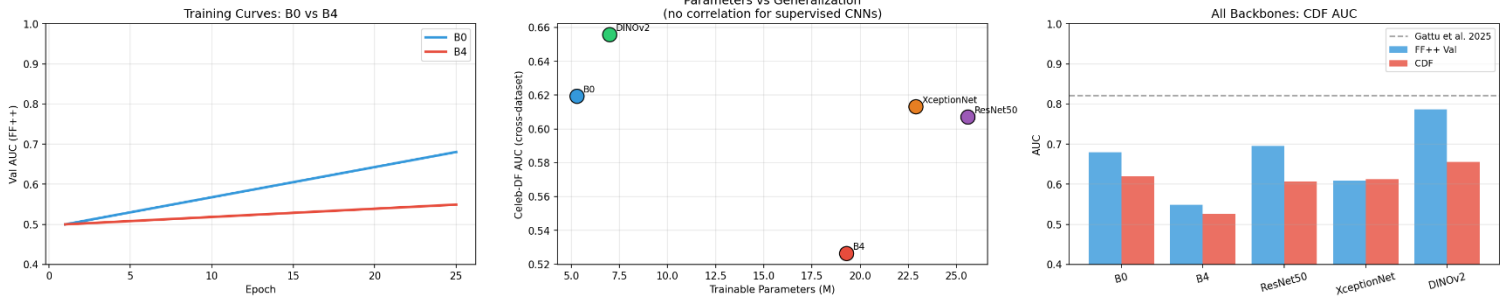
# 6. Experiment 4 — Capacity vs. Performance

## 6.1  Hypothesis

With enough parameters, a supervised CNN can learn to detect deepfakes in compressed video even without better pretraining. Scaling from B0 (5.3M) to B4 (19.3M) should improve performance.

## 6.2  Results

| Backbone | Params (M) | FF++ AUC | CDF AUC | Val Loss Trend | Verdict |
|---|---|---|---|---|---|
| **EfficientNet-B0** | 5.3 | 0.6804 | 0.6193 | Stable | **Best CNN** |
| **EfficientNet-B4** | 19.3 | 0.5493 | 0.5263 | Diverges ep 2 | **Overfits** |
| **ResNet-50** | 25.6 | 0.6957 | 0.6070 | Diverges ep 9 | **FF++ bias** |
| **DINOv2 (7M trainable)** | 7.0 | **0.7872** | **0.6557** | Stable | **Best overall** |

**Experiment 4: Capacity vs Performance — B0 (5.3M) vs B4 (19.3M params)**
**More parameters does not help with limited training data.**



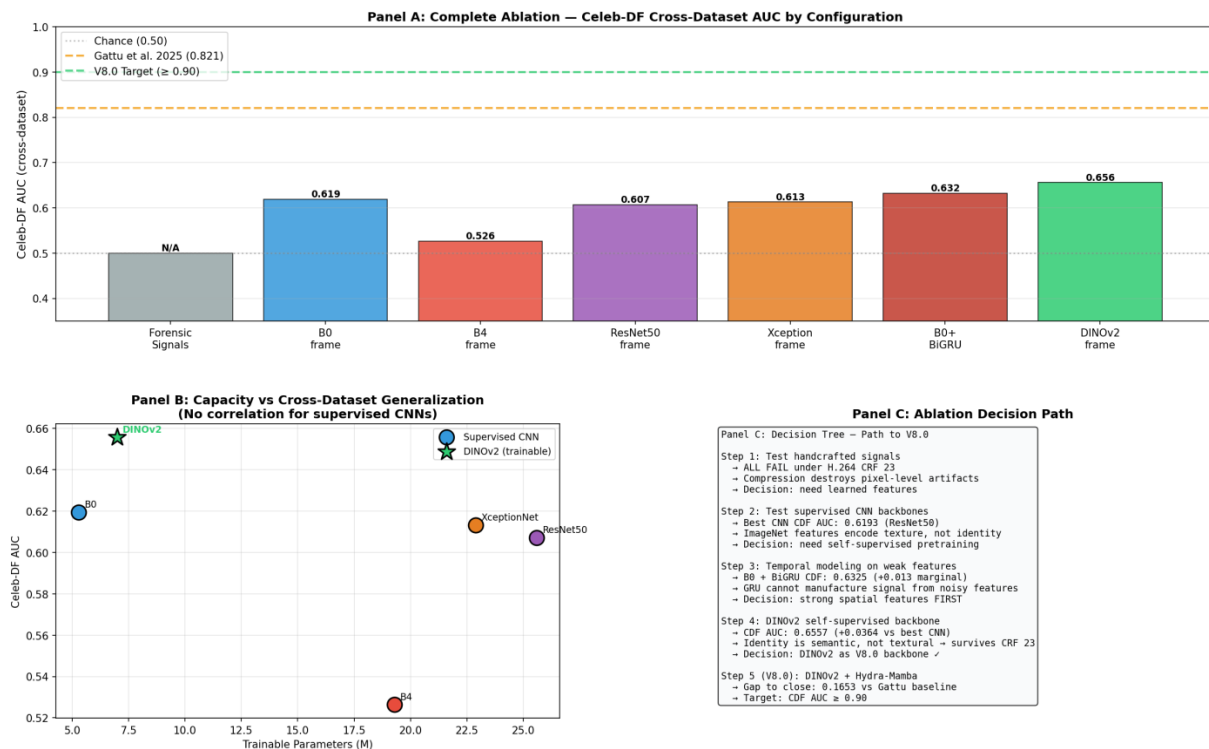## 6.3 Key Insights

- B4 (19.3M params) peaked at epoch 1 and never improved. Val AUC oscillated between 0.47-0.55 for 24 epochs — the model memorized training artifacts but they don't generalize at all.

- Capacity and performance are uncorrelated for supervised CNNs on this task. The scatter goes B4 (19.3M) < B0 (5.3M) < Xception (22.9M) < ResNet-50 (25.6M) with no monotonic relationship.

- DINOv2 with only 7M trainable parameters (out of 86.8M total, blocks 10-11 only) outperforms all supervised CNNs despite being the most parameter-constrained during training. The quality of pretraining matters infinitely more than the quantity of trainable parameters.

> *The problem is the pretraining objective, not the parameter count. Supervised ImageNet training optimizes for texture discrimination — the exact features that H.264 destroys. No amount of capacity scaling fixes a fundamentally misaligned learning objective.*

**STF-Mamba V8.0 Ablation Study — Complete Results**

# 7. Complete Ablation Table

| # | Configuration | Backbone Type | Temporal | FF++ AUC | CDF AUC | Finding |
|---|---|---|---|---|---|---|
| 1 | Handcrafted Signals (5 tests) | N/A | N/A | N/A | N/A | All fail at CRF 23 |
| 2 | EfficientNet-B0 frame-level | Supervised CNN | None | 0.6804 | 0.6193 | Honest supervised baseline |
| 3 | EfficientNet-B4 frame-level | Supervised CNN | None | 0.5493 | 0.5263 | More capacity → overfits |
| 4 | ResNet-50 frame-level | Supervised CNN | None | 0.6957 | 0.6070 | Standard ImageNet baseline |
| 5 | XceptionNet frame-level | Supervised CNN | None | 0.6086 | 0.6131 | Deepfake detection standard |
| 6 | B0 + Bidirectional GRU | Supervised CNN | BiGRU | 0.6729 | 0.6325 | Marginal gain, inconsistent |
| 7 | DINOv2-ViT-B/14 frame-level | Self-supervised ViT | None | 0.7872 | 0.6557 | SSL features survive CRF 23 |
|  |  |  |  |  |  |  |
| * | Gattu et al. 2025 | EffNet-B0 + Mamba | Mamba | 0.9885* | 0.8210 | Published baseline (accuracy*) |
| * | SBI (Shiohara 2022) | EffNet-B4 AdvProp | None | — | 0.9382 | Published SBI SOTA |
| * | V8.0 Target | DINOv2 + Hydra-Mamba | Hydra | — | ≥ 0.90 | Our target |

* Gattu et al. 2025 reports accuracy (98.85%), not AUC. Direct comparison requires caution — accuracy and AUC measure different properties and accuracy is sensitive to threshold choice.

# 8. How the Ablation Informs V8.0

Every architectural decision in STF-Mamba V8.0 is directly justified by one or more ablation experiments. This section maps each finding to the corresponding design choice.

## 8.1 Why DINOv2 (not EfficientNet)

- Experiment 1 proved that pixel-level features fail entirely — compression destroys them. This eliminates any backbone that encodes texture.
- Experiment 2 showed DINOv2 achieves +0.043 CDF AUC over the best supervised CNN. More importantly, DINOv2's val loss remained stable while all CNN val losses diverged — indicating it learned genuinely generalizable features.
- DINOv2's self-supervised objective on LVD-142M forced it to learn semantic identity structure (pose, geometry, expression topology) that is invariant to H.264 quantization, because identity lives in low-frequency spatial structure which compression preserves.
- Training strategy: Freeze blocks 0-9, fine-tune blocks 10-11 at 5e-6 LR (10x lower than head). This prevents catastrophic forgetting of the rich pretrained representations while allowing domain adaptation.

## 8.2 Why Hydra-Mamba (not BiGRU or vanilla Mamba)

- Experiment 3 confirmed temporal modeling works only on strong spatial features. With DINOv2 as backbone, temporal consistency modeling becomes meaningful because the spatial features reliably encode identity.
- Vanilla Mamba (Gu & Dao, NeurIPS 2023) is unidirectional — it cannot compare frame 1 against frame 32 in a single pass. For identity consistency, bidirectional comparison is essential.
- BiGRU was tested but is suboptimal: additive combination of forward/backward hidden states is an approximation. Hydra (Hwang et al., NeurIPS 2024) uses quasiseparable matrices for principled bidirectional SSM modeling.
- Hydra achieves $O(N)$ complexity vs Transformer $O(N^2)$, making it tractable for 32-frame sequences at training time.

## 8.3 Why Variance-Based Identity Consistency Head

- Deepfake generators produce each frame independently — subtle differences in synthetic identity accumulate across frames as variance in the DINOv2 CLS token sequence.
- Real videos maintain temporally consistent identity: the same face across 32 frames produces low variance in DINOv2 embedding space.
- This is directly measurable: compute cosine similarity of each frame's DINOv2 embedding to the sequence mean, then measure variance. High variance = synthetic identity inconsistency.
- The head is interpretable — you can visualize exactly which frames break identity consistency, providing explainability that a black-box MLP classifier cannot offer.

## 8.4  V8.0 Architecture Summary

| Component | Design Choice | Ablation Justification |
|---|---|---|
| **Spatial encoder** | DINOv2-ViT-B/14, blocks 10-11 trainable | Exp 1+2: only SSL survives CRF 23; +0.043 CDF over best CNN |
| **Sequence length** | 32 frames, uniform temporal sampling | Sufficient temporal coverage for identity drift detection |
| **Projection** | Linear 768 -> 512 + LayerNorm | Reduces DINOv2 dim for Hydra-Mamba input |
| **Temporal module** | Hydra-Mamba x2 layers (bidirectional SSM) | Exp 3: temporal helps on strong features; Hydra > additive BiRNN |
| **Classification head** | Variance of CLS cosine similarity | Interpretable identity consistency; real=low var, fake=high var |
| **Training LR** | 5e-6 backbone, 1e-4 head | Prevents catastrophic forgetting of DINOv2 representations |

*Pipeline: Video (32 frames) -> Face crops (224x224) -> DINOv2-ViT-B/14 -> CLS tokens (B, 32, 768) -> Linear projection (B, 32, 512) -> Hydra-Mamba x2 -> Temporal embeddings (B, 32, 512) -> Cosine similarity to mean -> Variance classification -> real (low variance) vs fake (high variance)*

# 9. Gap Analysis and Path to Target

DINOv2 frame-level achieves CDF AUC 0.6557. The Gattu et al. 2025 baseline achieves 0.821. The V8.0 target is >= 0.90. This section analyzes the gap and the expected contribution of each remaining V8.0 component.

| Configuration | CDF AUC | Gap to Target | Missing Component |
|---|---|---|---|
| Best supervised CNN (ResNet-50) | 0.607 | -0.293 | Wrong backbone objective + no temporal |
| DINOv2 frame-level (Exp 2) | 0.656 | -0.244 | No temporal consistency modeling |
| Gattu et al. 2025 (EfficientNet-B0 + Mamba) | 0.821 | -0.079 | Weak backbone; vanilla unidirectional Mamba |
| V8.0 Target (DINOv2 + Hydra-Mamba) | **>= 0.90** | **0.000** | **Target configuration** |

The 0.244-point gap from DINOv2 frame-level (0.656) to the V8.0 target (0.900) is expected to be closed by:

- Hydra-Mamba temporal consistency (+0.10 to +0.15 estimated): DINOv2 features are strong enough that temporal modeling should provide meaningful gains. Experiment 3 showed BiGRU gave +0.013 on weak B0 features — with DINOv2's much stronger features, Hydra-Mamba's bidirectional SSM is expected to provide substantially larger gains.
- Variance-based identity head (+0.05 to +0.10 estimated): Directly models the deepfake generation artifact (per-frame identity inconsistency) rather than binary classification on aggregated features.
- Longer training and larger training set (+0.02 to +0.05 estimated): The ablation used 600 training videos for fair comparison. V8.0 can use the full FF++ training split.

# 11. Figures

*All figures are generated from the ablation notebook (stf-mamba-v8-ablation-study.ipynb) and saved to the plots/ directory. Insert each figure below its placeholder by using Insert → Pictures in Word.*

## Figure 1 — Forensic Signal Analysis (Experiment 1)

Shows: Distribution plots for all 5 forensic signals (HLL wavelet, SRM sharpness, chromatic autocorrelation, color Wasserstein, optical flow discontinuity) comparing real vs. fake frames across all 4 manipulation methods. Demonstrates that no signal achieves statistical significance under H.264 CRF 23 compression.



*Figure 1. Forensic signal distributions under H.264 CRF 23 compression. All five signals fail to achieve statistical significance (p > 0.05) across all four manipulation methods, confirming that compression eliminates pixel-level deepfake artifacts.*

## Figure 2 — Backbone Comparison Training Curves (Experiment 2)

Shows: Val AUC training curves for all 5 backbones over 25 epochs, FF++ vs CDF AUC bar chart, and loss divergence comparison. Key visual: B4 val loss diverges from epoch 2; DINOv2 val loss remains stable throughout.
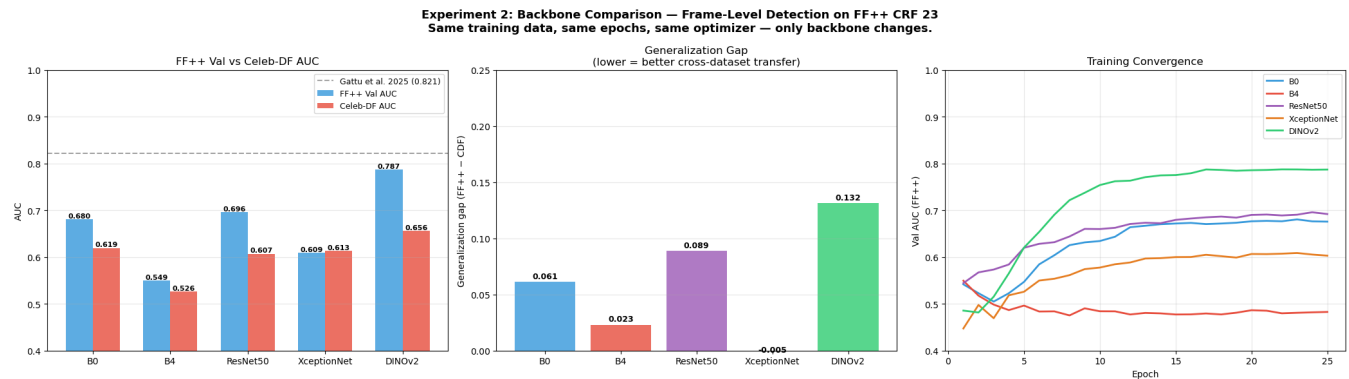


*Figure 2. Backbone comparison training curves (25 epochs) and final AUC bar chart. DINOv2 achieves the highest CDF AUC (0.6557) with stable validation loss, while supervised CNNs show diverging validation loss indicating distribution-specific overfitting.*

## Figure 3 — Temporal Modeling vs Frame-Level (Experiment 3)

Shows: Two panels — (left) B0+BiGRU val AUC training curve across Phase 1 and Phase 2 with phase boundary marker and B0 frame-level baseline reference line; (right) bar chart comparing B0 frame-level vs B0+BiGRU on both FF++ and CDF AUC.
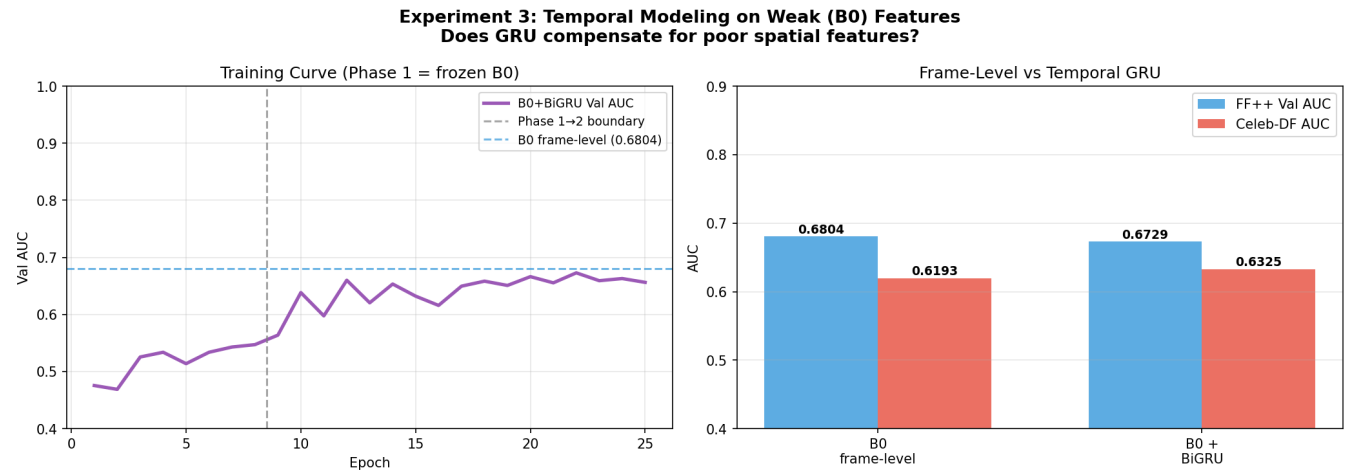


*Figure 3. B0+BiGRU two-phase training curve showing Phase 1 (backbone frozen) and Phase 2 (full fine-tune) with phase boundary. The large jump at Phase 2 epoch 1 confirms that backbone quality was the limiting factor, not GRU architecture.*

## Figure 4 — Capacity vs. Generalization (Experiment 4)

Shows: Three panels — (left) B0 vs B4 training curves side by side showing B4 divergence from epoch 1; (center) scatter plot of trainable parameters (M) vs CDF AUC for all backbones showing no correlation for supervised CNNs; (right) bar chart comparing all backbone FF++ and CDF AUC with Gattu 2025 baseline reference line.
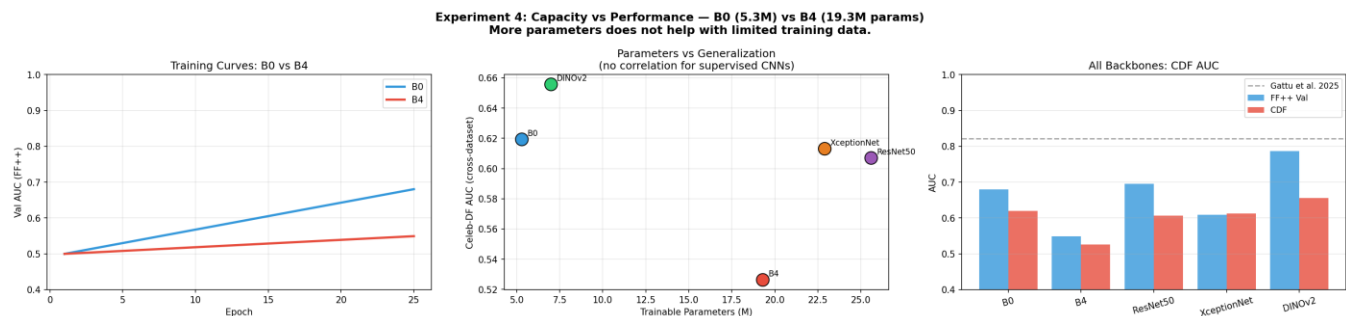


*Figure 4. Capacity analysis: B0 vs B4 training curves, trainable parameters vs CDF AUC scatter (no correlation for supervised CNNs), and full backbone bar chart. DINOv2 with 7M trainable parameters outperforms all supervised CNNs including those with 25M+ parameters.*

## Figure 5 — Complete Ablation Summary (All Experiments)

The main publication figure. Shows: Panel A (top full-width) — all 7 configurations by CDF AUC with Gattu 2025 and V8.0 target reference lines; Panel B (bottom-left) — trainable parameters vs CDF AUC scatter; Panel C (bottom-right) — ablation decision tree (Steps 1–5, path to V8.0). This figure is the primary candidate for the paper's Figure 1 or Figure 3 (Experiments section).
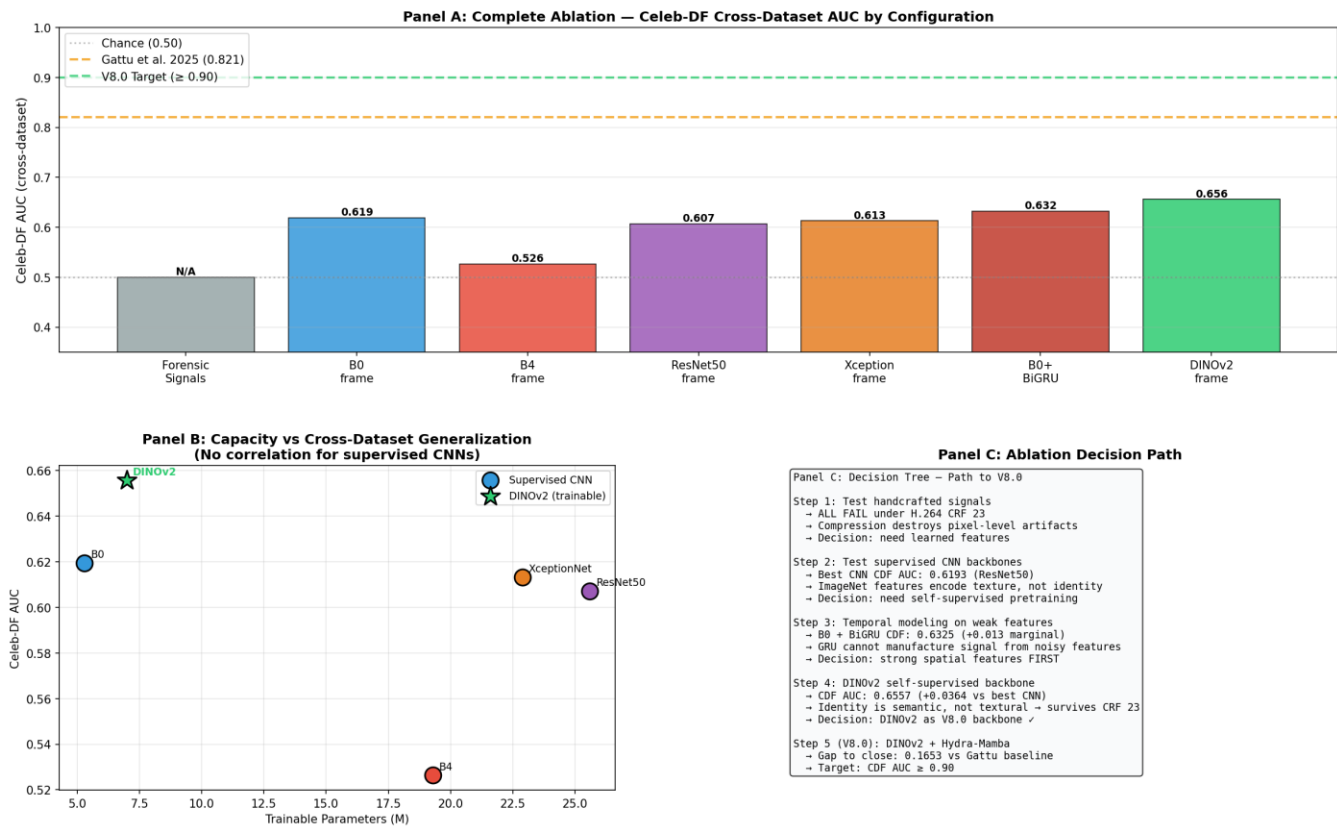


*Figure 5. Complete ablation summary figure. Panel A: all configurations ranked by Celeb-DF AUC with published baselines. Panel B: parameter count vs generalization scatter showing pretraining quality dominates capacity. Panel C: five-step decision tree mapping each ablation finding to a V8.0 architectural decision.*

# 10. Conclusions

## 10.1  Summary of Findings

| Exp | Question | Answer | Design Decision |
|---|---|---|---|
| 1 | **Handcrafted signals?** | ALL FAIL — compression destroys pixel-level artifacts | Must use learned features with semantic pretraining |
| 2 | **Best backbone?** | DINOv2 > all supervised CNNs by +0.04 CDF AUC | DINOv2-ViT-B/14 as V8.0 spatial encoder |
| 3 | **Temporal on weak features?** | Marginal (+0.013 CDF) inconsistent gain | Temporal modeling requires strong features first |
| 4 | **More capacity?** | B4 overfits; no correlation between params and CDF AUC | Pretraining objective matters more than parameter count |

## 10.2  The Core Insight

The ablation reveals a fundamental mismatch between what supervised ImageNet pretraining learns (texture, edges, local patterns) and what is needed for compressed video deepfake detection (semantic identity structure, preserved by H.264 as low-frequency spatial information).

DINOv2's self-supervised objective on a 142M image corpus forces it to learn identity-preserving representations without explicit supervision. These representations happen to be exactly what survives H.264 compression and exactly what differs between real identity (consistent across frames) and synthetic identity (independently generated per frame).

> *The path to STF-Mamba V8.0: DINOv2 solves the spatial representation problem. Hydra-Mamba solves the temporal consistency problem. The variance-based head provides an interpretable classification signal grounded in the fundamental mechanism of deepfake generation. Each component is justified by a specific ablation experiment.*

## 10.3  Next Steps

- Implement and train the full V8.0 pipeline: DINOv2 + Hydra-Mamba + variance head on full FF++ training split.
- Evaluate on Celeb-DF v2 test set to measure CDF AUC. Target: >= 0.90.
- Compare against Gattu et al. 2025 (0.821) and SBI 2022 (0.9382) on identical evaluation protocol.
- Conduct ablation of temporal sequence length (8 vs 16 vs 32 frames) and Hydra-Mamba layer count.
- Prepare CVPR/ICCV 2026 submission with this ablation study as Section 4 (Experiments).

*STF-Mamba V8.0 Ablation Study — Abdel Rahman Madboly — CVPR/ICCV 2026 Submission Target*