Redacted author name and title investigates whether stochastic inference (single-pass and Monte Carlo dropout), temperature scaling, and ensembles can improve the reliability of political deepfake detectors beyond point predictions. Using two CNN backbones (ResNet-18, EfficientNet-B4) on a politically filtered real–synthetic image dataset, the study evaluates calibration, proper scoring rules, and uncertainty–error alignment, including confidence-band analyses, OOD generator-disjoint tests, and robustness to JPEG compression. The paper concludes that uncertainty serves as a conditional, decision-level signal—most informative in high-confidence regimes for selective abstention and triage—while discrimination (AUC) remains largely unchanged across inference procedures.

| Category | Score | Reason |
|----------|-------|--------|
| Abstract | 9 | Clear, self-contained statement of objectives, methods, datasets, evaluation metrics, and main findings; readers can understand the study without the manuscript. |
| Recency | 9 | References are largely from 2023–2025 with foundational works; highly up-to-date for uncertainty and deepfake detection. |
| Scope | 8 | Covers conditional uncertainty-aware detection, political focus, CNN backbones, calibration, and OOD per the title and keywords; scope notes limitations appropriately. |
| Relevance | 9 | Addresses a high-stakes, timely problem (political deepfakes) with operationally relevant evaluation; avoids unnecessary background digressions. |

| | | |
|---|---|---|
| Factual Errors | 9 | No substantive factual errors detected; a few minor formula/notation typos do not affect correctness or conclusions. |
| Language | 8 | Generally precise, technical prose; minor typographic/spacing artifacts and occasional stylistic inconsistencies but scientifically sound tone. |
| Formatting | 8 | Manuscript is structured well with clear sections, tables, and metrics; some equation rendering artifacts and symbol formatting inconsistencies. |
| Suggestions | 8 | Introduces a clear, decision-oriented reliability framing and confidence-band uncertainty analysis; could add conformal prediction baselines, identity-disjoint splits, and more modern backbones (e.g., ViT) to broaden impact. |
| Problems | 8 | Targets gaps in calibration and uncertainty–error alignment for political deepfakes; highlights when uncertainty adds value beyond confidence and quantifies practical effect sizes; cautions against overinterpreting global AUROC of uncertainty. |
| Assumptions | 7 | Assumptions (e.g., metadata-based political filtering, matched-generator ID setup) are explicit and tested via ablations; external validity to unseen identities and platforms remains limited. |
| Consistency | 9 | Quantitative and qualitative findings cohere with literature (calibration "`discrimination; ensembles help; OOD harms calibration/accuracy); claims are properly scoped. |
| Robustness | 7 | Includes generator-disjoint OOD, JPEG robustness, MC T-sensitivity, and dropout-rate ablations; robustness to identity shift, heavy platform pipelines, or adversarial perturbations is untested. |
| Logic | 9 | Conclusions follow from the reported data and sensitivity analyses, carefully distinguishing discrimination from reliability and conditional utility of uncertainty. |
| Statistical Analysis | 9 | Appropriate use of ROC-AUC with DeLong tests (paired OOD), bootstrap CIs for accuracy/ECE/Brier, and proper scoring rules; clear threshold semantics and uncertainty–error AUROC calculations. |
| Controls | N/A | Wet-lab style experimental controls are not applicable; computational controls (baselines, ablations, fixed seeds) are provided elsewhere. |
| Corrections | 7 | Addresses influential factors via stratification (generators), compression sweeps, and preprocessing ablations; further correction for identity overlap or platform artifacts would strengthen claims. |

| | | |
|---|---|---|
| Range | 8 | Explores meaningful parameter ranges (MC T" {1,5,10,20,50}, dropout p" {0,0.1,0.2,0.5}, JPEG Q sweep, multiple normalizations/resolutions) capturing practical regimes. |
| Collinearity | N/A | Multicollinearity diagnostics are not relevant to this non-regression, deep learning classification setting. |
| Dimensional Analysis | 8 | Equations involve probabilities and losses (dimensionless) and are consistent; minor typographical artifacts do not affect dimensional correctness. |
| Experimental Design | 8 | Clear protocols, fixed seeds, hardware/software specs, and held-out validation/test; ensembles and MC are well controlled. Improvements: add identity-disjoint OOD, conformal prediction selective-risk curves, platform-specific degradations, and transformer backbones for completeness. |
| Ethical Standards | informational | Recommend adding an ethics statement covering the use of political images, data licenses, potential harms (false positives/negatives in political contexts), and deployment safeguards (human review, escalation policies, misuse prevention). |
| Conflict Of Interest | informational | Include an explicit conflict-of-interest and funding statement; if none, state: "The author declares no competing interests and no external funding." |
| Normalization | informational | Primary pipeline deviates from ImageNet normalization and later ablates normalization choices; retain this transparency and consider making the dataset-specific mean–std and applied transforms easily reproducible via a config file and data card. |
| Idea Incubator | informational | Cross-disciplinary analogies (neutral, heuristic mappings): - Economics (option pricing under volatility): Treat uncertainty as implied volatility; high-confidence/high-uncertainty cases are like options with deep-in-the-money deltas but elevated vega, guiding selective abstention like hedging exposure. - Epidemiology (test sensitivity/specificity vs. prevalence): As base rate (deepfake prevalence) shifts, calibrated probabilities and abstentions resemble targeted screening thresholds; uncertainty flags subpopulations with higher false-positive risk. - Control theory (robust MPC): Confidence is nominal model prediction; uncertainty is model–plant mismatch estimate; selective abstention acts as a constraint tightening in high-risk states to keep the closed-loop system safe. - Information theory (rate–distortion): Coverage is rate; error under abstention is distortion; uncertainty-conditioned rejection traces a risk–coverage curve akin to allocating bits where marginal information gain is highest. - Ecology (predator–prey with refuges): High-uncertainty predictions act as refuges where the detector (predator) abstains, preserving system stability by preventing overconfident misclassifications that could cascade misinformation dynamics. - Queueing systems (triage under load): Uncertainty guides priority routing to human review queues; risk-aware thresholds minimize expected cost given finite reviewer capacity and arrival variability. |

| | | |
|---|---|---|
| Improve Citability | informational | To maximize reuse and citations: (1) Publish a data card detailing political keyword lists, generator identities, licenses, and known biases. (2) Release train/val/test and generator-disjoint OOD splits with checksums. (3) Provide a one-command reproducibility script and environment lockfile. (4) Add conformal prediction and risk–coverage API to enable downstream selective-prediction studies. (5) Include pre-registered evaluation protocols for future comparability. (6) Provide per-example prediction files (scores, calibrated scores, entropies, errors) for meta-analyses. (7) Add identity-disjoint and platform-degradation benchmarks. (8) Offer model zoos (weights) for each backbone and inference mode with versioned configs. |