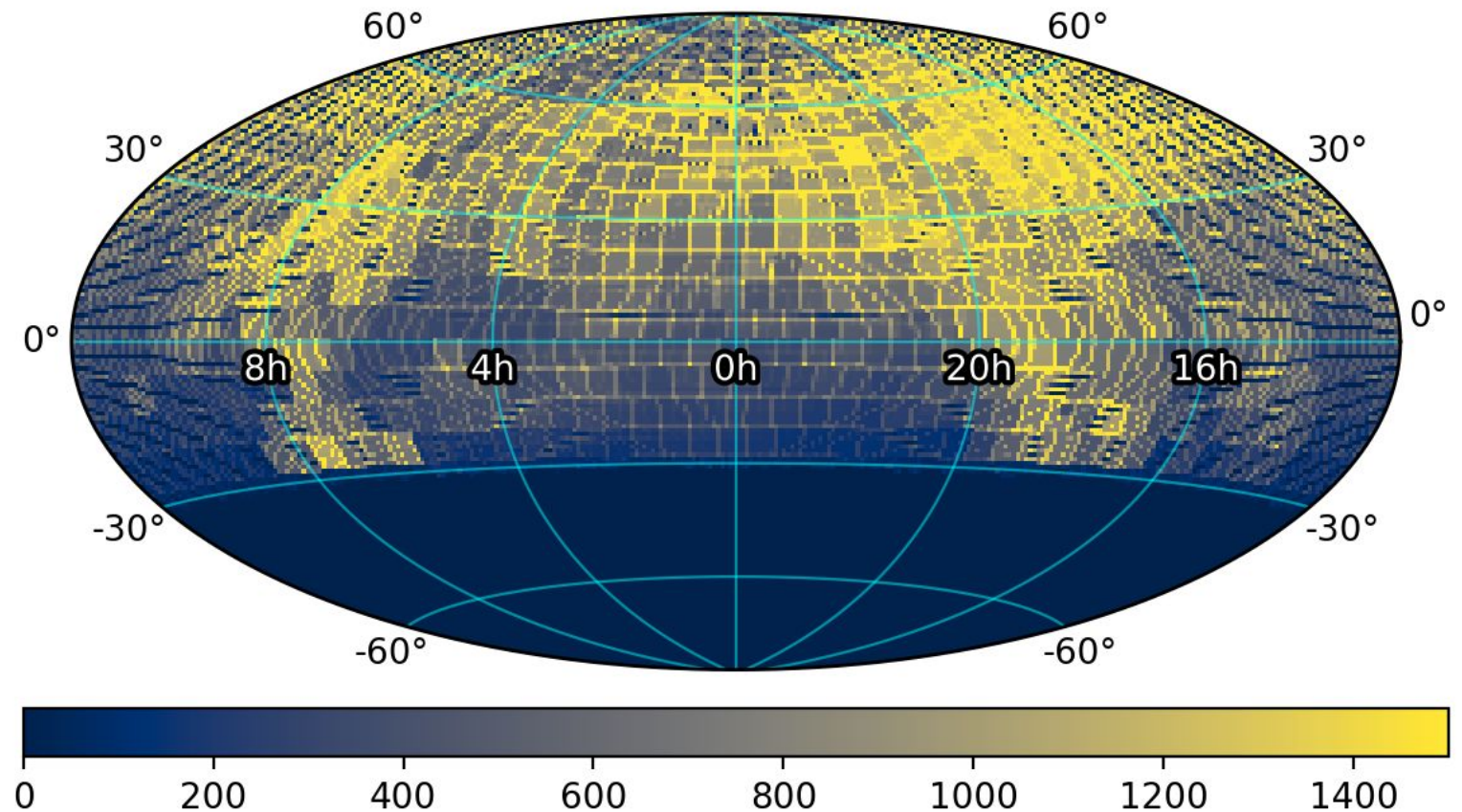


Deep Learning Applications in Astronomy



AY 119 · 2026

Ashish Mahabal — Caltech
Surveys, pipelines, decisions

From theory to practice

19 papers, multiple domains, mostly Caltech-connected.

Recurring theme: representation tricks matter more than fancy architectures.

The data deluge

ZTF: ~1M alerts per night → LSST/Roman coming

TESS: 200,000+ stars at 2-min cadence.

LIGO / Virgo / KAGRA: continuous strain, glitch-dominated.

CHIME / DSA-110: thousands of FRB candidates per day.

There is no human-only path forward. DL is the production layer.

Roadmap

1. Images — real/bogus, asteroids, subtraction-free
2. Light curves — DMDT, SCOPE, TESS
3. Spectra — SNIascore, CCSNscore
4. Gravitational waves — DMDT in LIGO + the GWSkyNet family
5. Radio transients — Frabjous (FRBs)
6. Language models — AstroAlertBench

1. Images

ZTF and the alert stream

Difference imaging: subtract a reference image, look for what's new.

Most candidates are artifacts: CR hits, optical ghosts, bad subtractions.

Allows a fast, high-recall gate before science pipelines spend time.

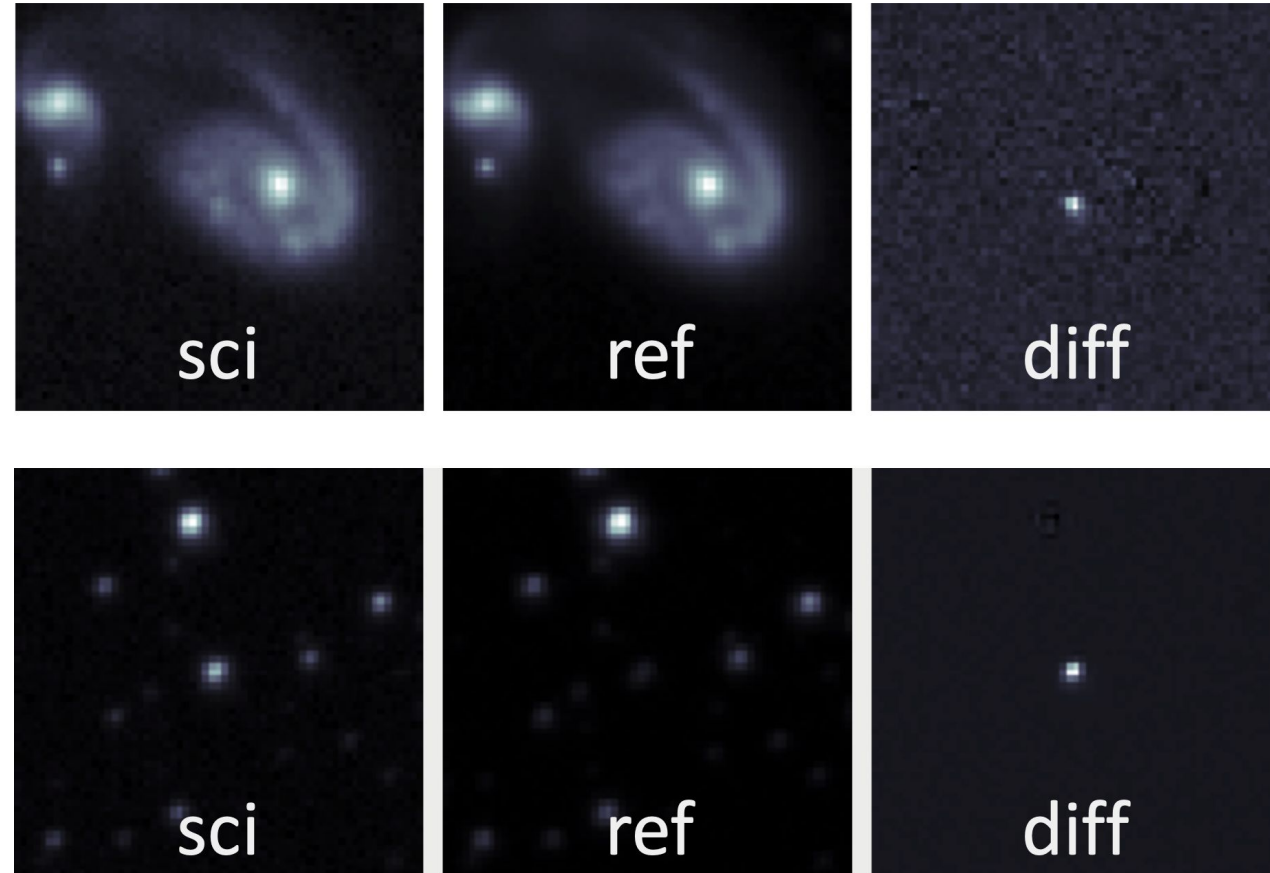
BRAAI — real or bogus (Duev+ 2019)

VGG6 CNN on (science, reference, difference) image triplets.

Trained on ~32k vetted ZTF examples; <20 min on one GPU.

At threshold 0.5: 1.1%
false-negative + 2.9%
false-positive on test.

Every ZTF science pipeline sits
behind this gate.



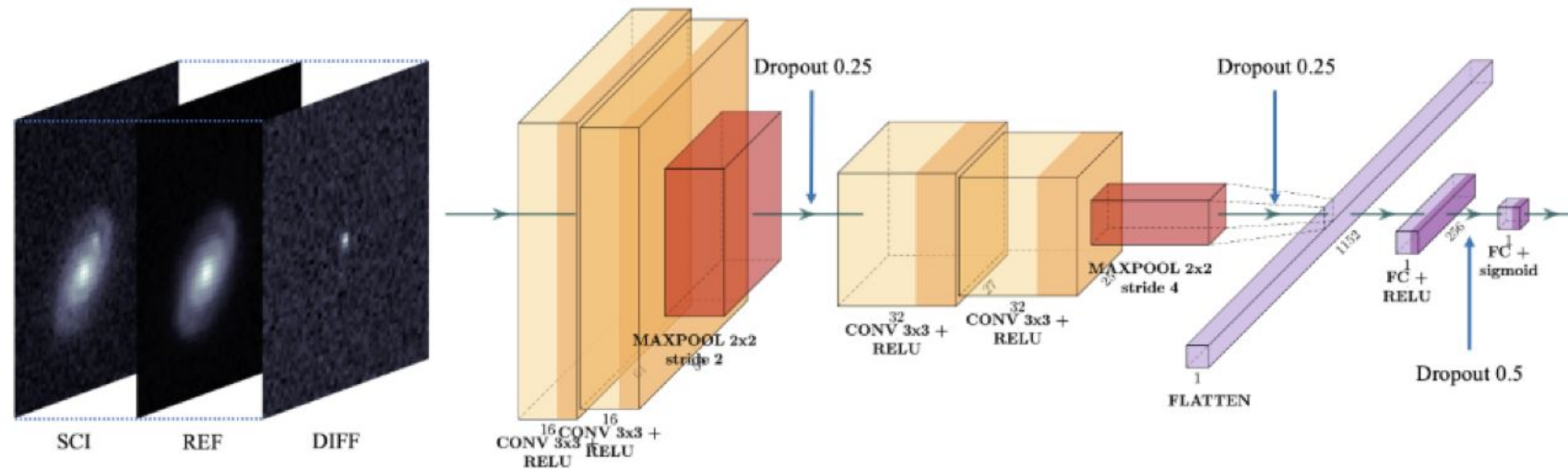
ZTF Real-Bogus separation ('braai')



Sci

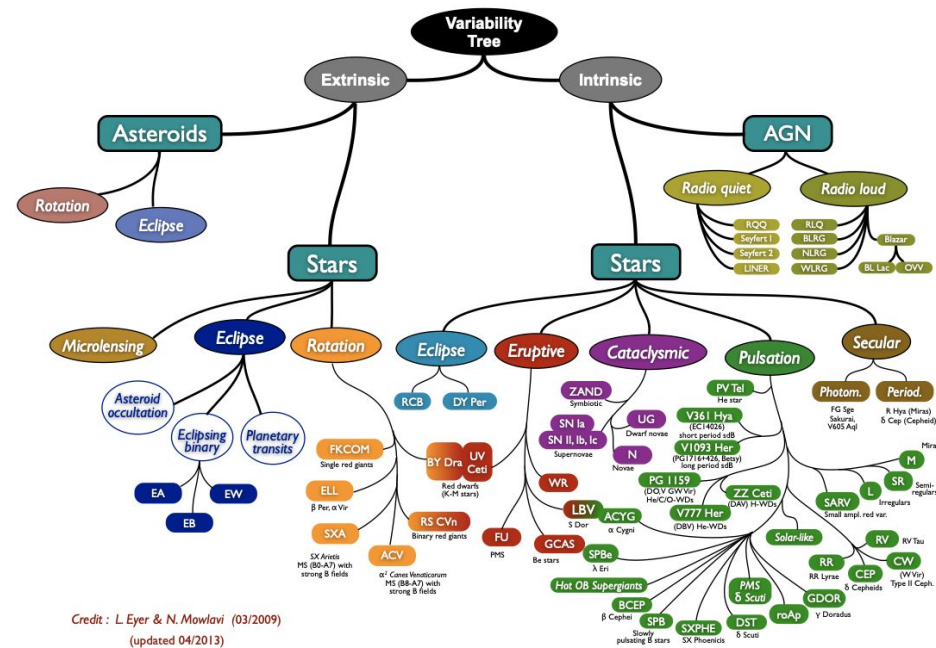
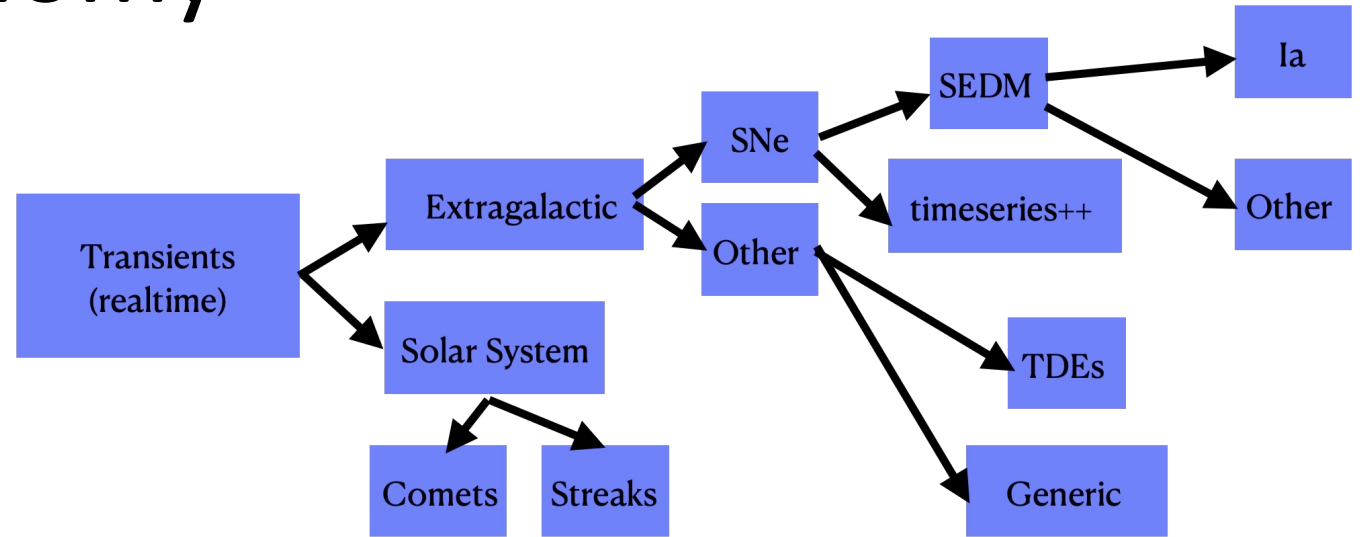
Ref

Diff



By area of astronomy

- Transients
- Detection
- Classification
- Variable stars
- Solar System



Credit : L. Eyer & N. Mowlavi (03/2009)
(updated 04/2013)

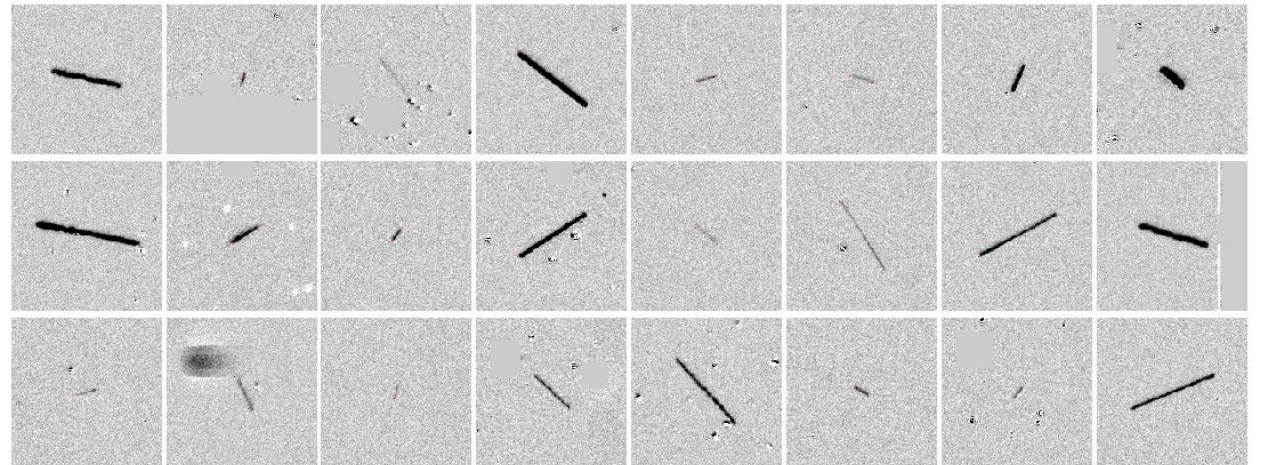
DeepStreaks — fast-moving asteroids (Duev+ 2019)

30-second exposures leave streaks for fast NEOs.

Trained on 1k real + 8k synthetic streaks + 6k bogus images.

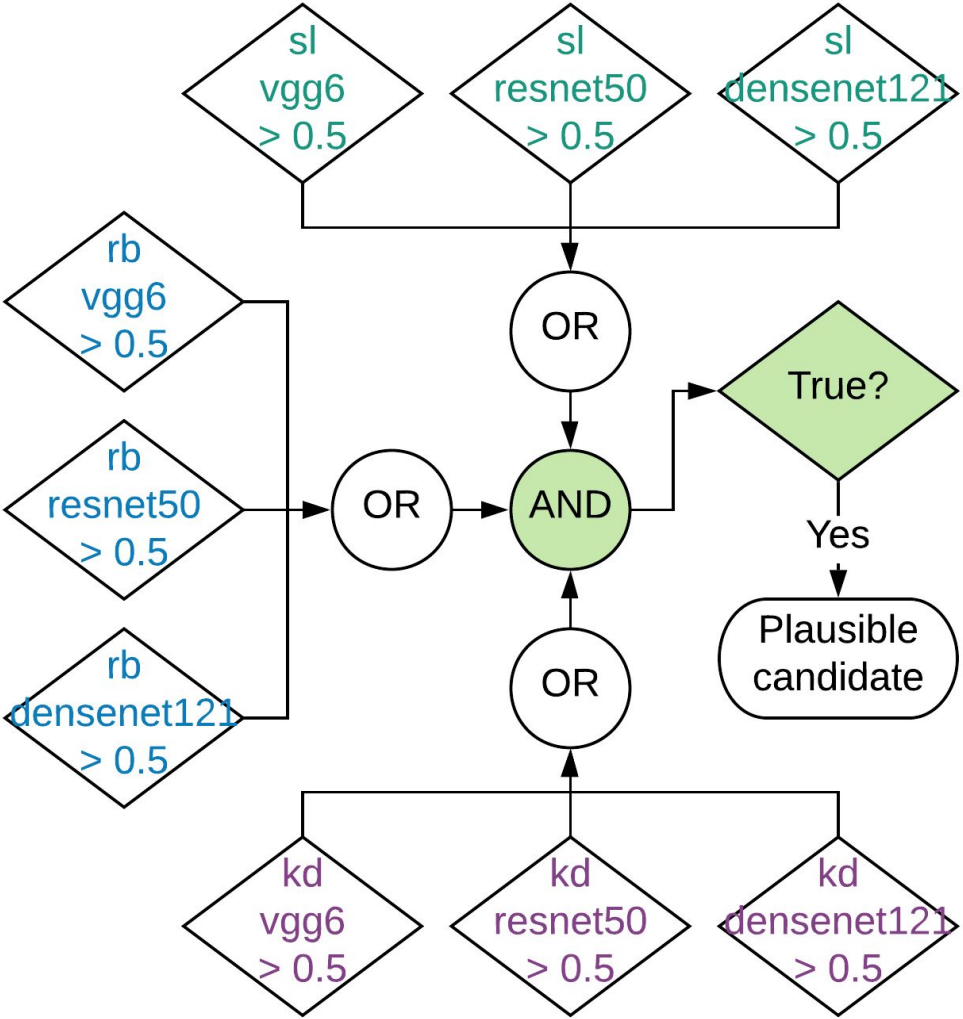
Ensemble: VGG6 + ResNet50 + DenseNet121, per task.

96–98% TPR; 50× fewer false positives than the RF baseline.



Real ZTF streaks — [arXiv:1904.05920](https://arxiv.org/abs/1904.05920)

DeepStreaks — pipeline architecture



TransiNet — subtraction-free detection (Sedaghat+ 2017)

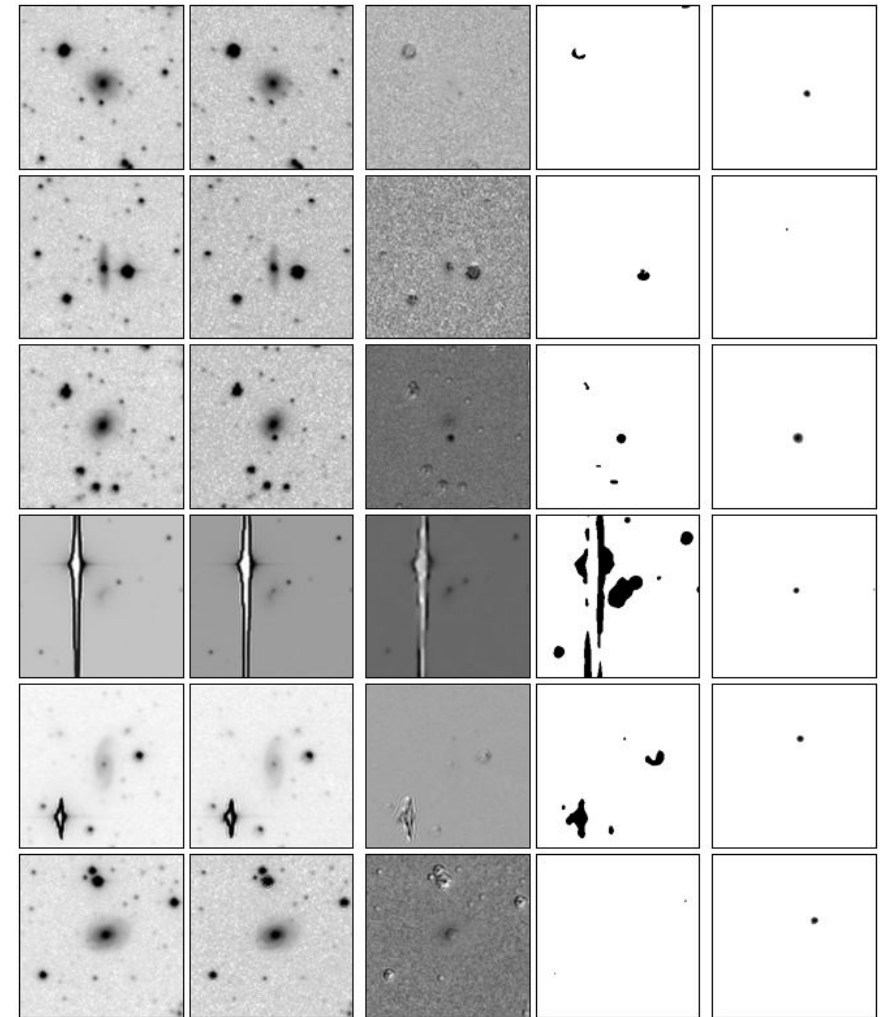
Skip the classical reference-subtraction pipeline entirely.

Encoder–decoder CNN: raw image pair \rightarrow clean detection map.

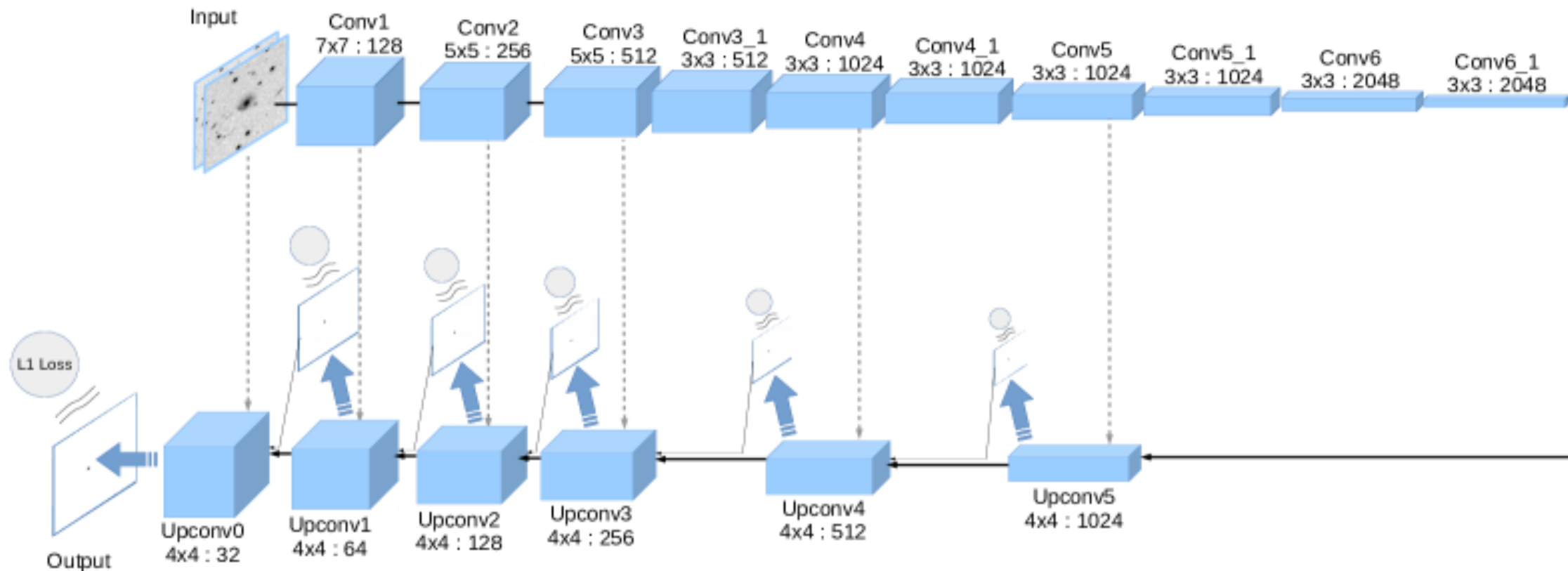
Trained on ~ 100 CRTS SNe + Galaxy Zoo synthetics; L1 loss at all scales.

75.5% recall at 98.4% precision on real SNe.

Hint of what foundation models could do for whole pipelines.



TransiNet — encoder-decoder architecture



2. Light curves

The light curve representation problem

Irregular sampling, gaps, heteroscedastic errors.

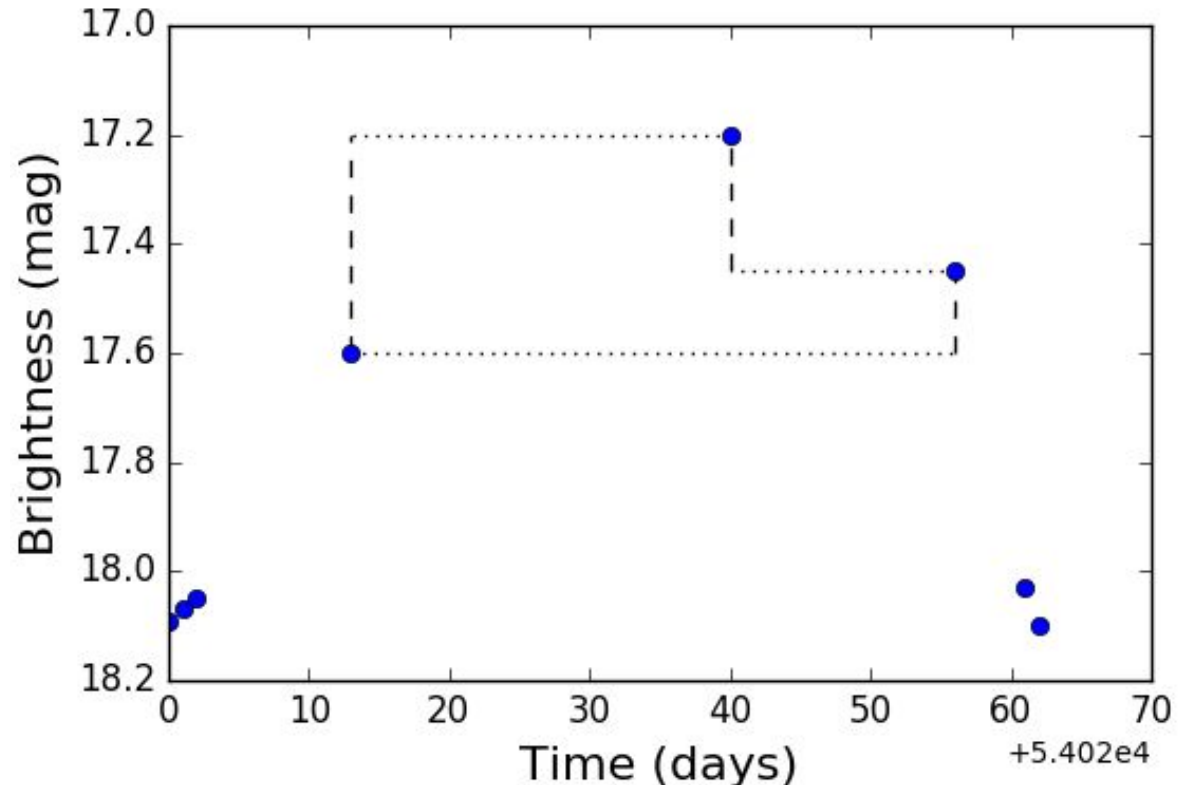
Can't feed naïvely into a CNN.

Three options: hand-crafted features, RNNs / Transformers, or re-package as images.

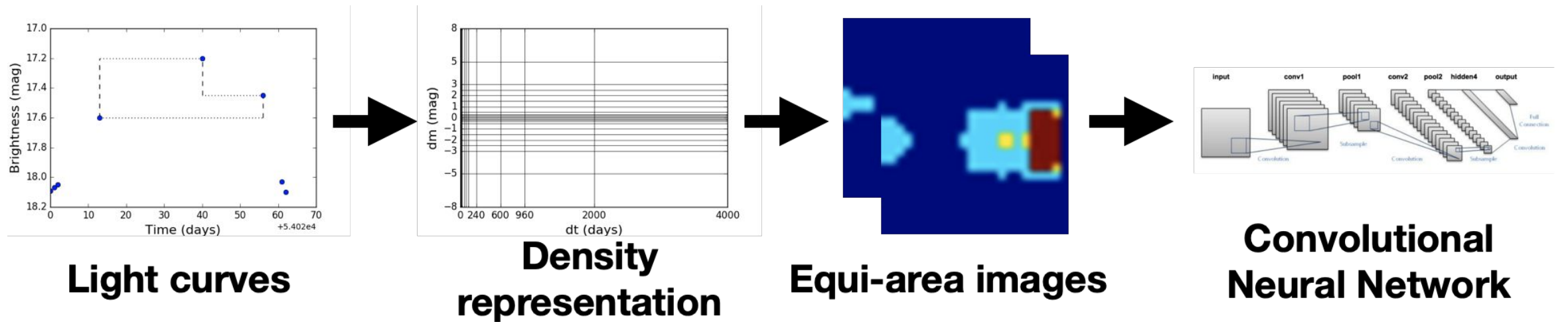
Our path: re-package as images — and then a CNN.

DMDT images (Mahabal+ 2017)

- For every pair of points compute $(\Delta\text{mag}, \Delta t)$.
- Bin into a 23×24 grid \rightarrow a grayscale image.
- Now a standard CNN classifies variable star types.
- 83% accuracy on 7-class CRTS; shallow CNN \approx deep CNN.
- Representation $>$ architecture.



DMDT — stacked class composites



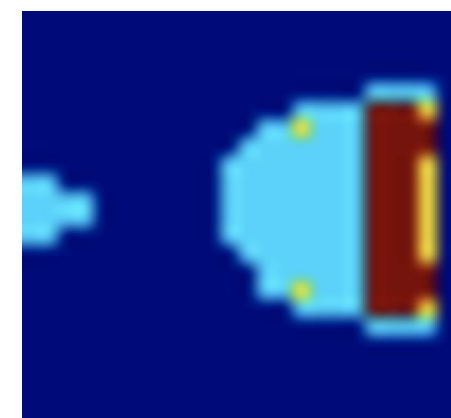
dmdt
images



RR Lyrae



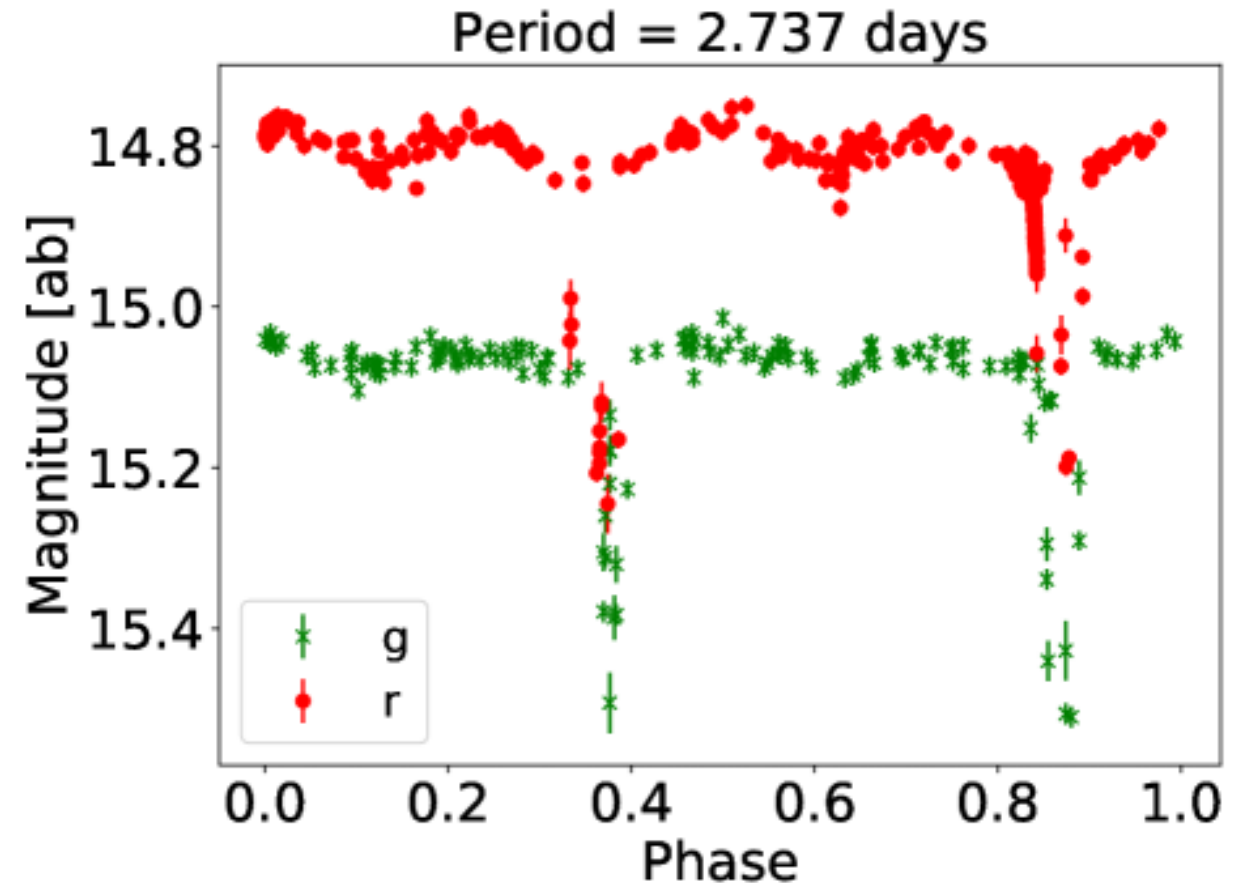
RS CVn



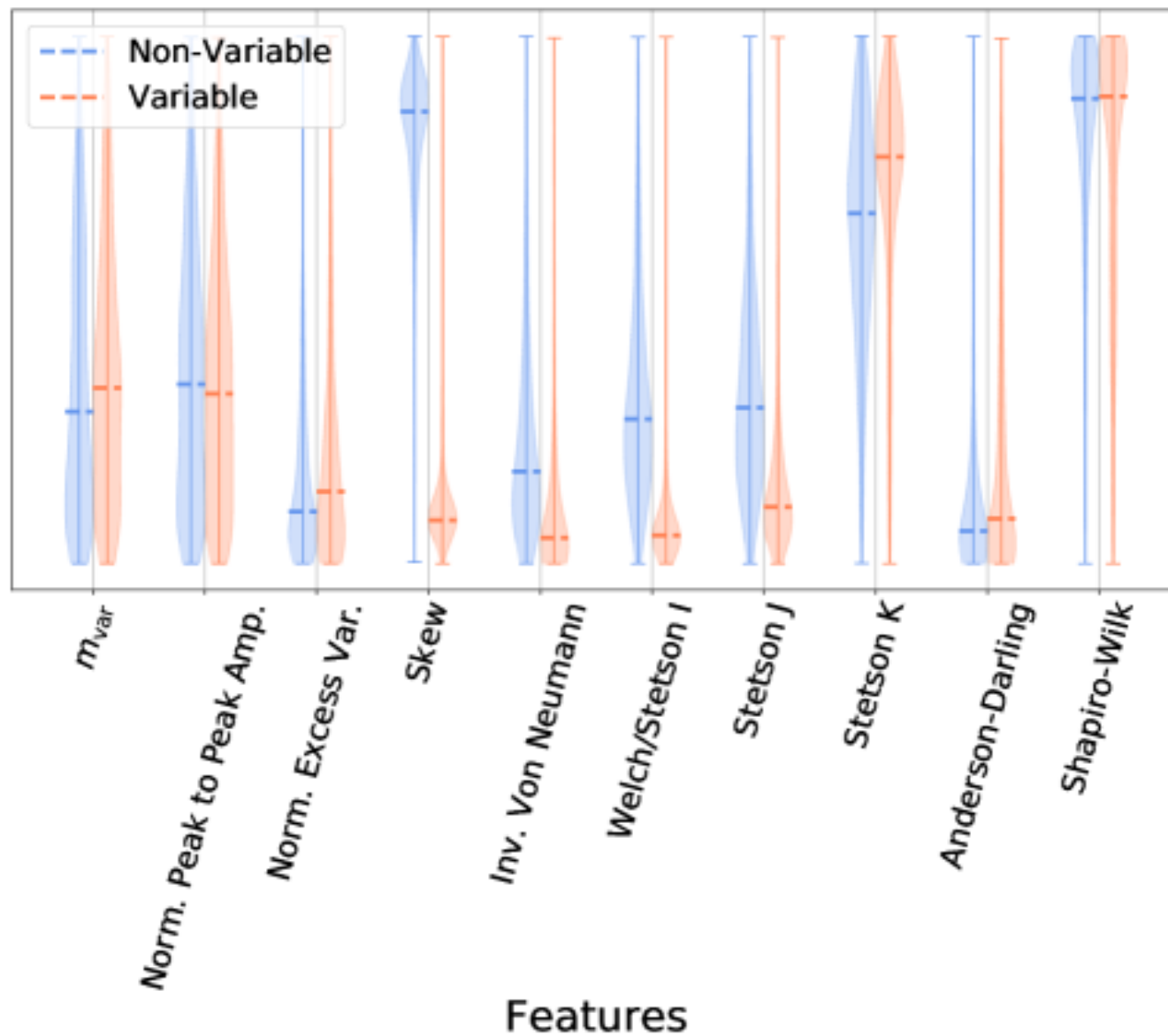
LPV

ztfperiodic — period-finding at scale (Coughlin+ 2020)

- Before classification: variability stats + a period.
- GPU Conditional Entropy + Lomb–Scargle, run in parallel.
- 37 statistical features per light curve, aliasing mitigation.
- ~1 second / light curve → all 1.3B ZTF sources reachable.
- The plumbing under SCOPE.



ztfperiodic — GPU scaling



SCOPE — Source Classification Project (van Roestel+ 2021)

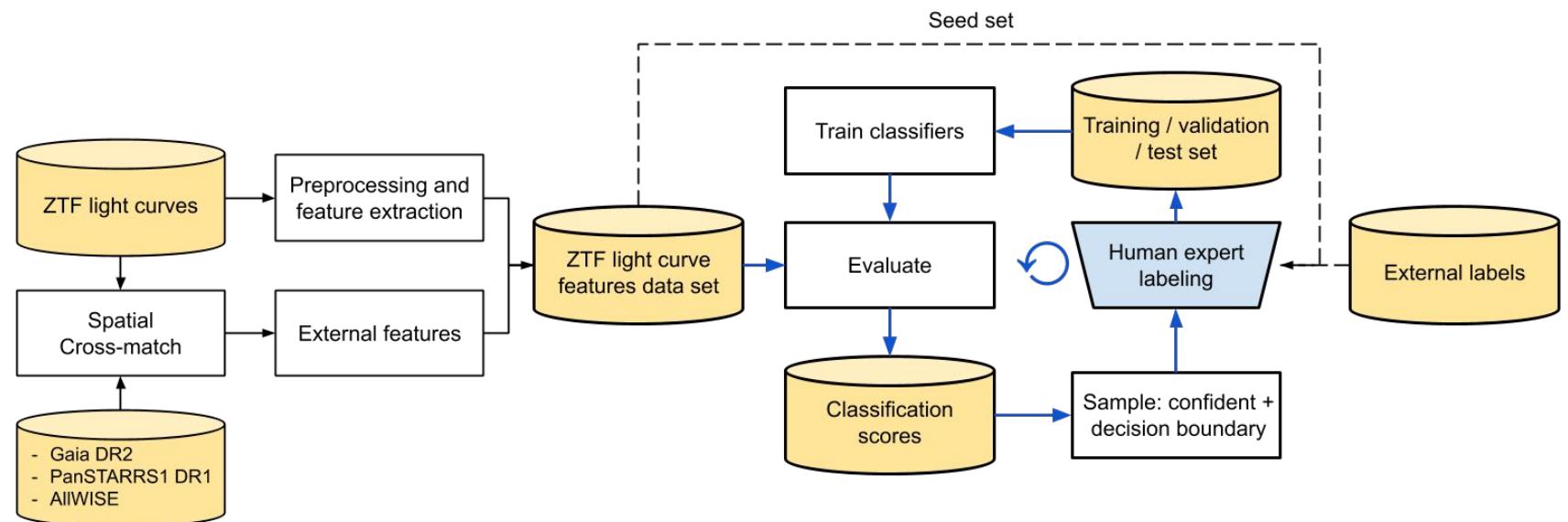
Many dichotomous binary classifiers, not one big multi-class head.

Two taxonomies in parallel: ontological + phenomenological.

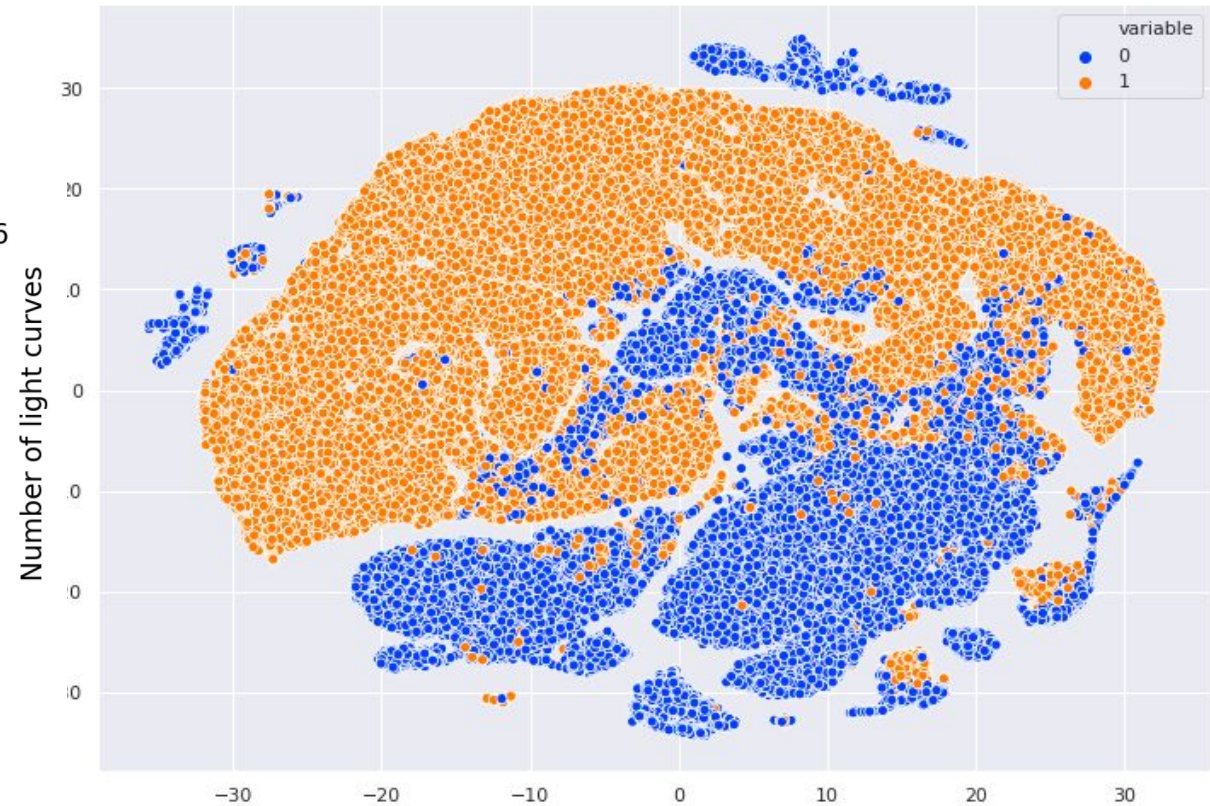
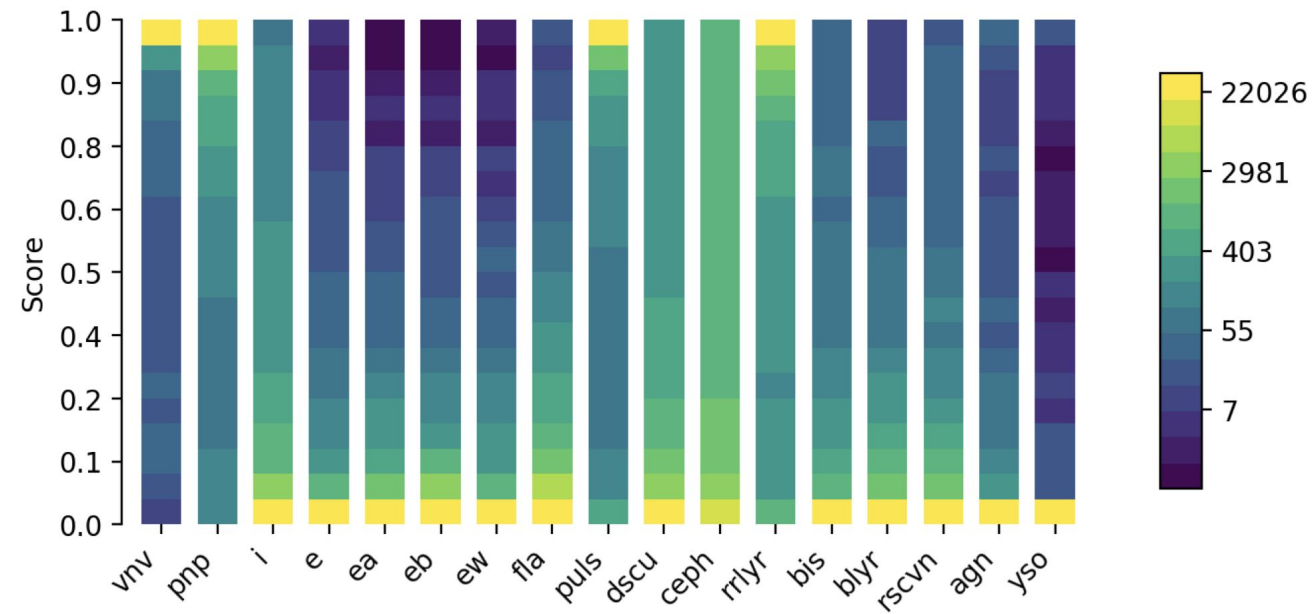
40 ZTF features + DMDT + Gaia/Pan-STARRS/AllWISE crossmatches.

DNN and XGBoost both trained; active-learning loop via SkyPortal.

89% precision for RR Lyrae extraction in production.



SCOPE — feature-space structure



SCOPE at full ZTF scale (Healy+ 2023)

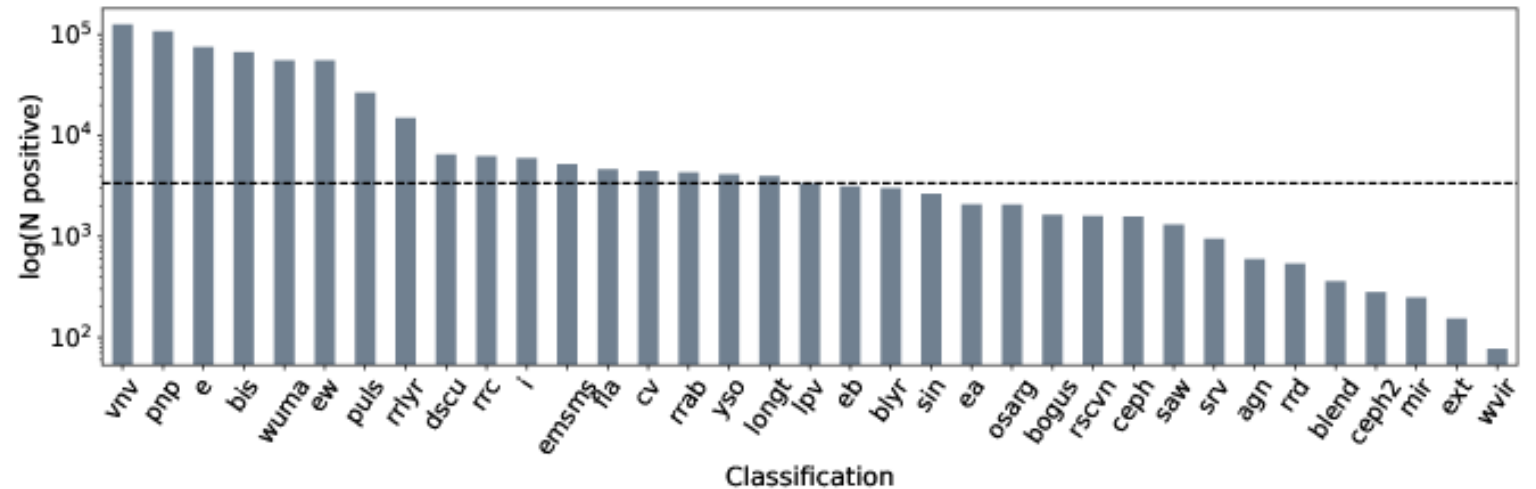
Same dichotomous framework, now applied to all of ZTF.

35+ classifiers; 170k hand-labelled light curves after iterative AL.

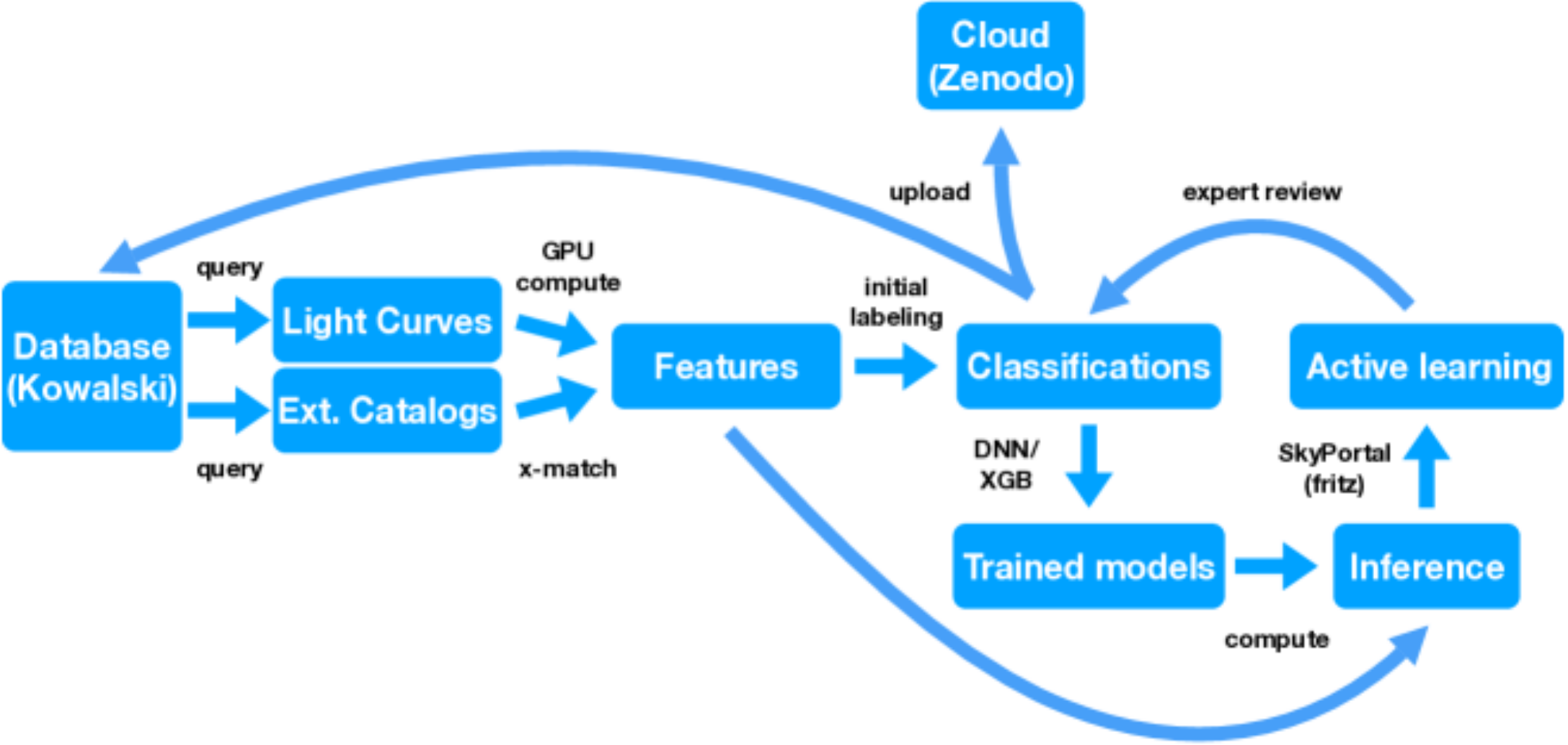
112,476,749 classifications produced (40 fields).

Period determination + Fourier power dominate feature importance.

Public catalog — the deliverable for the community.



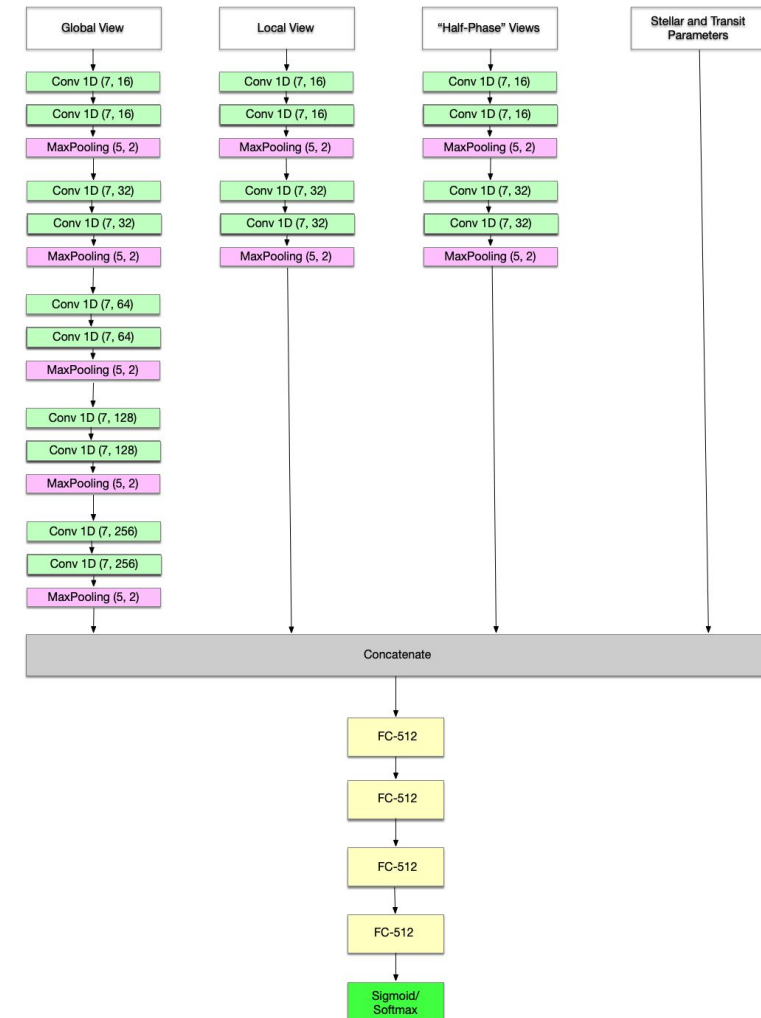
SCOPE — end-to-end workflow



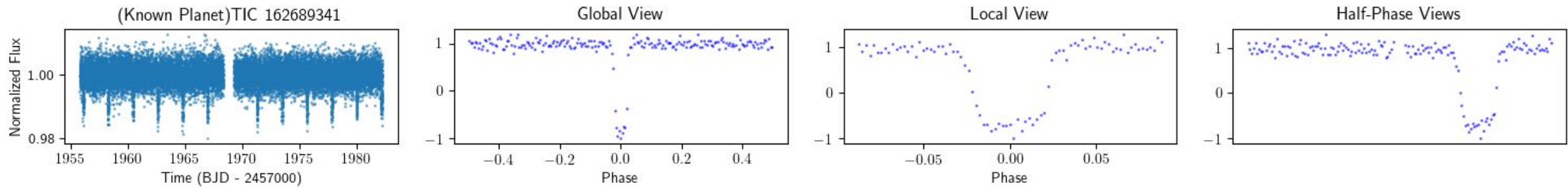
From raw light curves through active learning to public catalog — arXiv:2312.00143

TESS exoplanet transits (Rao+ 2021)

- 2-min cadence light curves; Transit-Least-Squares pre-screen.
- Three input views into one CNN: global (201 bins) + local (81) + half-phase.
- Half-phase view kills eclipsing-binary contaminants.
- AUC 0.908; recovers 91.9% of known planets at score > 0.7 .
- 38 new vetted candidates across 7 sectors.



TESS — phase-folded example



3. Spectra

SNiascore — real-time Type Ia typing (Fremling+ 2021)

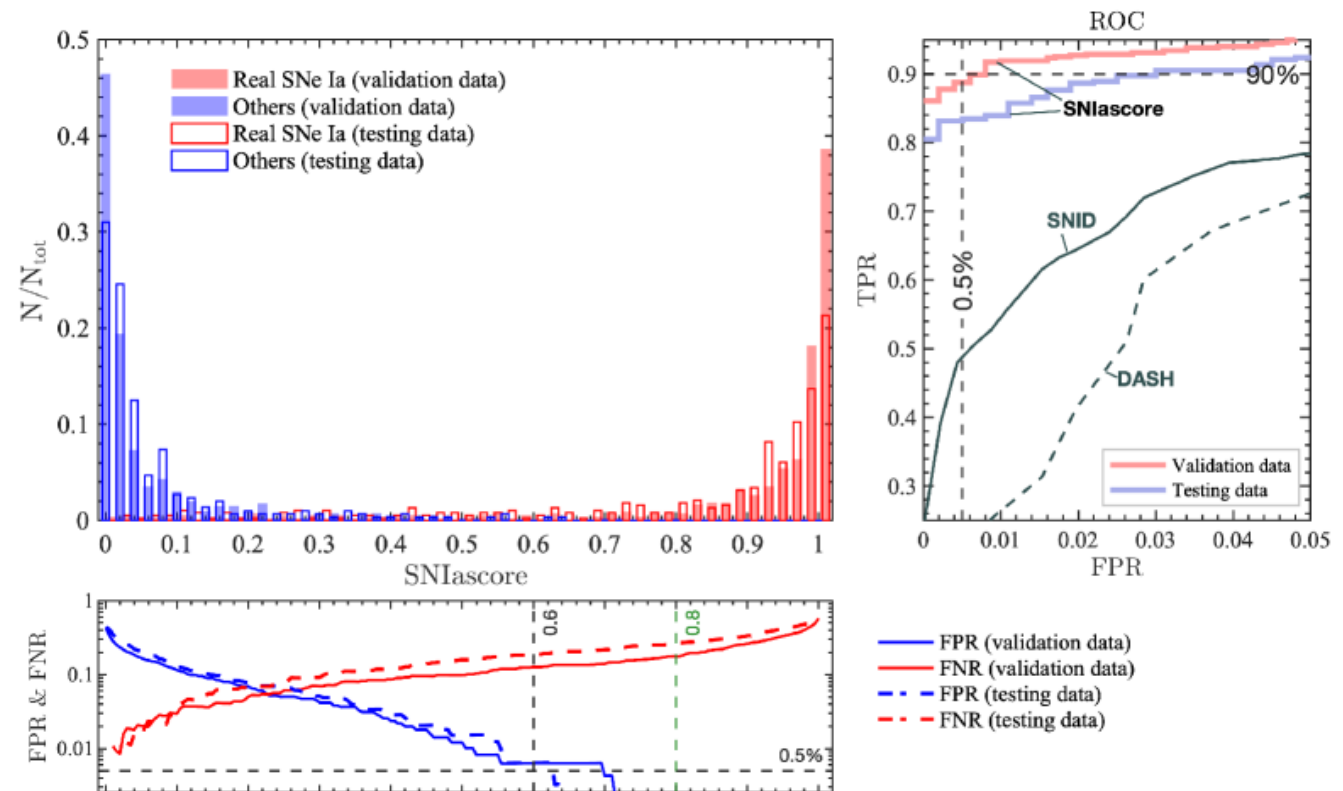
Low-resolution ($R \sim 100$) SEDM spectra \rightarrow SN Ia / not-Ia.

BiLSTM + GRU + BiLSTM with 40–45% dropout.

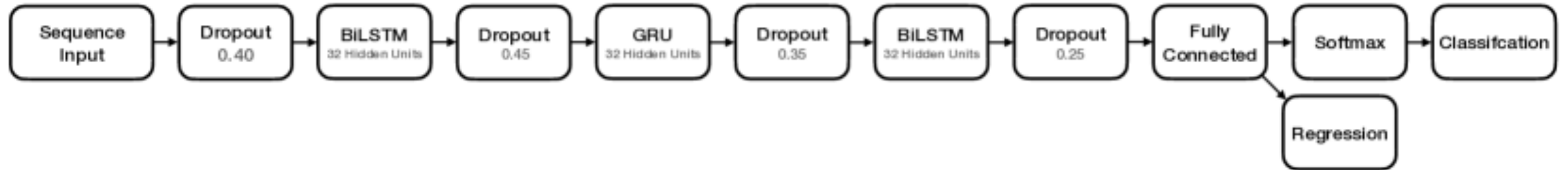
Joint redshift regression head (custom MBE loss).

90% TPR at 0.6% FPR on validation; runs inside the live BTS workflow.

Frees humans for the rare / weird objects.



SNIscore — RNN architecture



CCSNscore — core-collapse SNe (Sharma+ 2024)

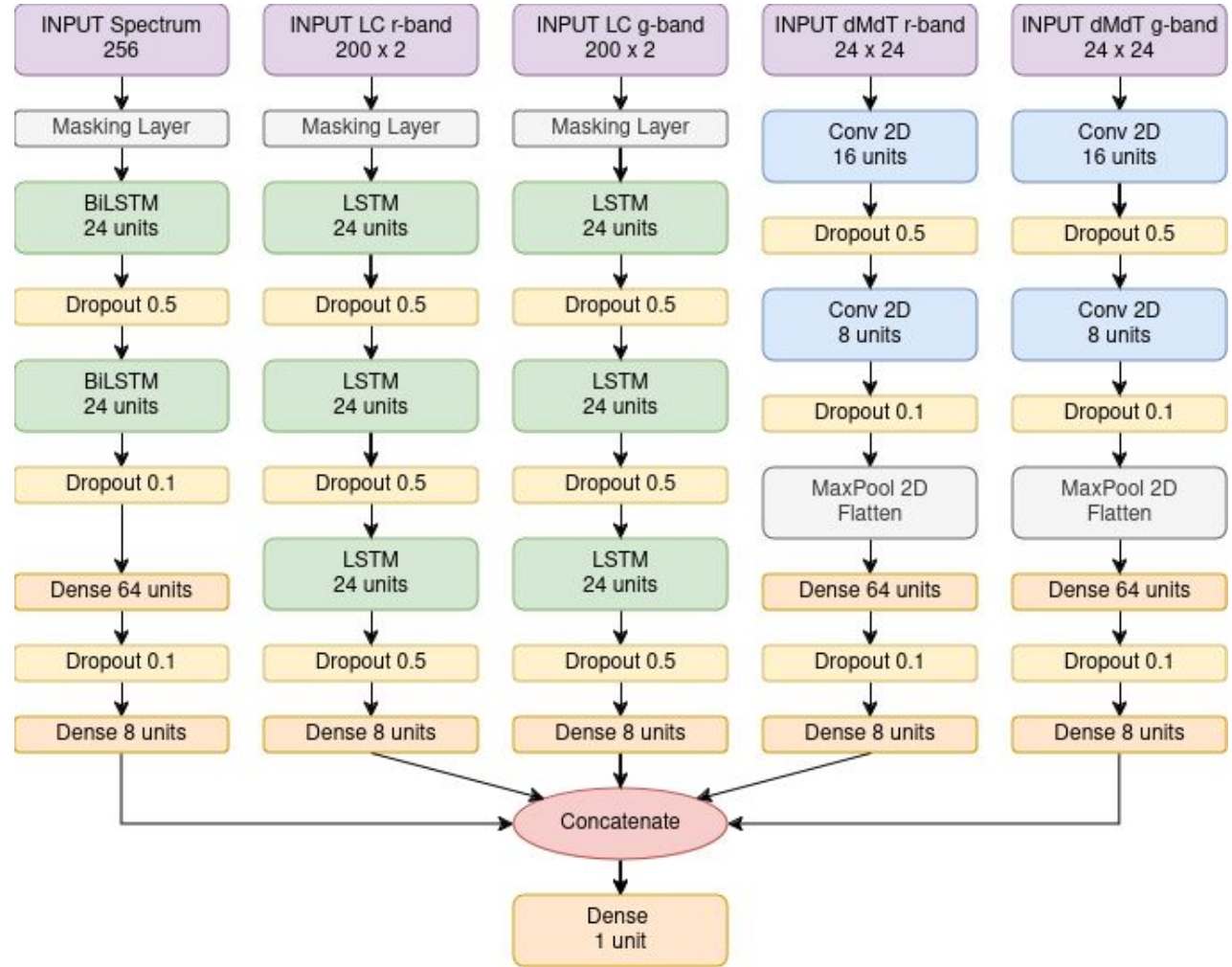
Multi-input: spectra + ZTF g/r light curves + DMDT.

Hierarchical: H-rich vs. H-poor, then II / IIb / IIc or Ib / Ic / Ic-BL.

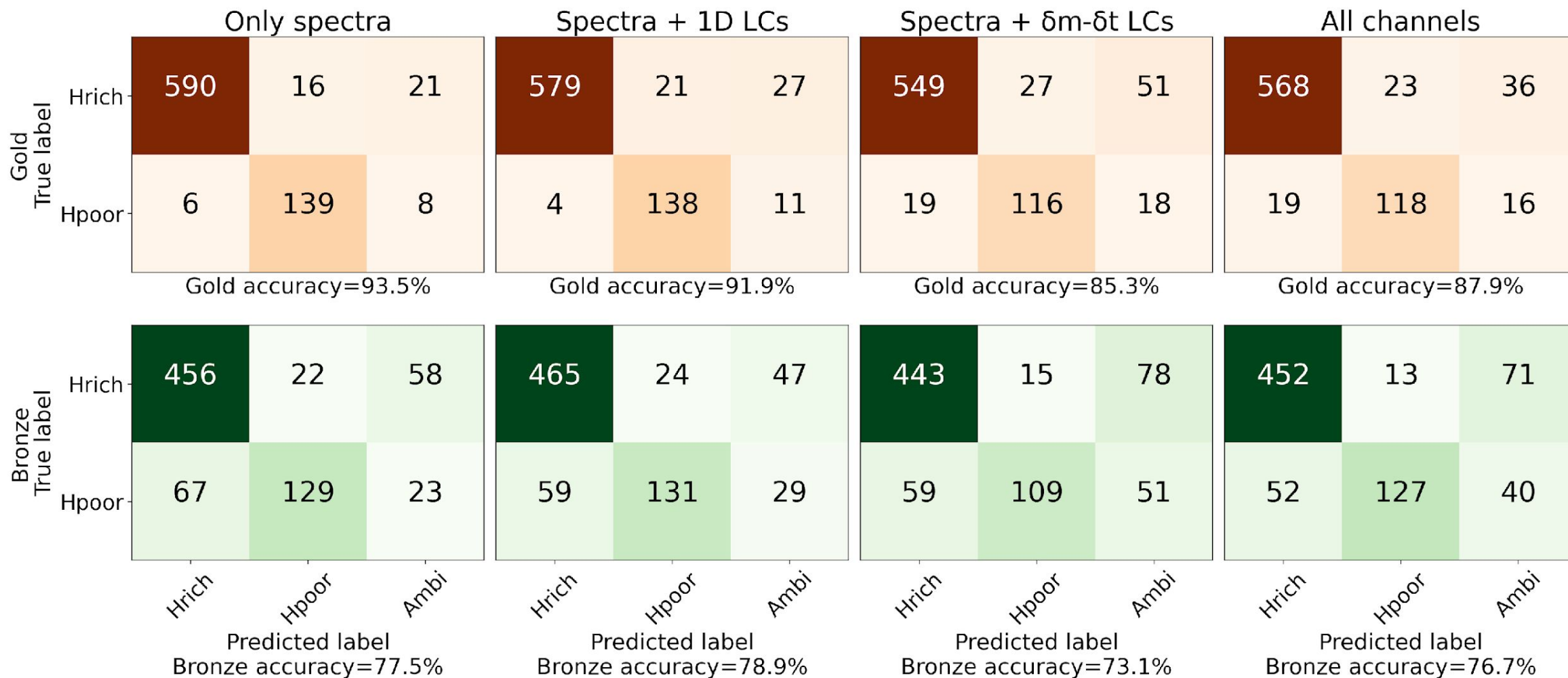
Bi-LSTM on spectra, LSTM on photometry.

MC dropout ($\times 100$) \rightarrow per-prediction uncertainty.

98% accuracy on gold-quality spectra.



CCSNscore — Layer-2 confusion matrices



4. Gravitational waves

GW data 101

Strain time series, kHz sampling, multiple detectors.

Q-transform \rightarrow 2-D spectrogram (now an image, CNN-friendly).

Noise is dominated by transient glitches — the real adversary.

DL tasks: glitch ID, real/bogus, source classification, parameter inference, and... predicting lock loss.

DMDT goes to LIGO (Biswas+ 2019)

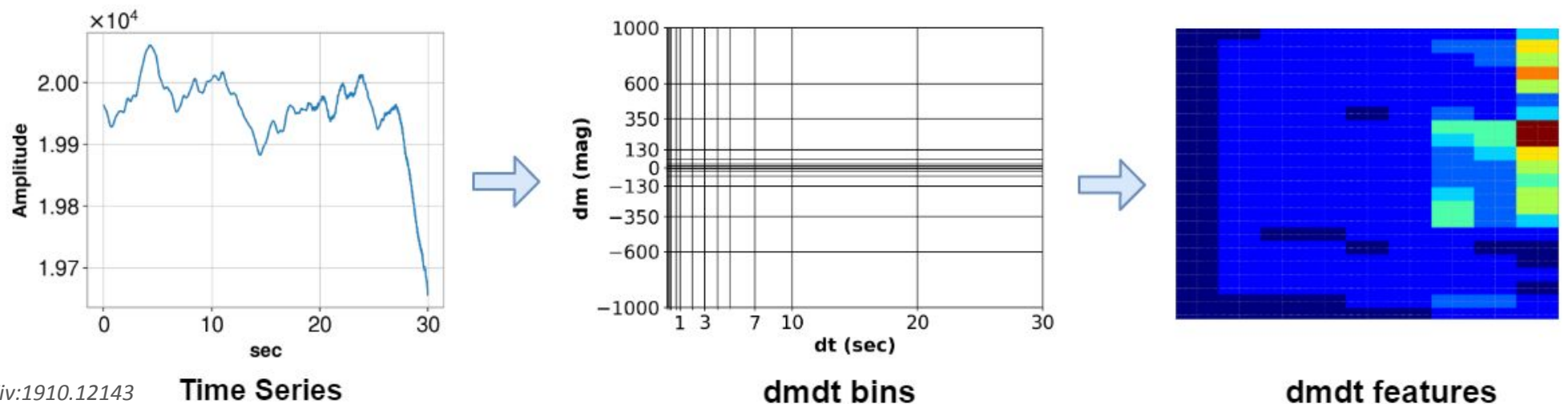
The DMDT representation applied to detector control channels.

22 auxiliary channels \rightarrow DMDT images \rightarrow CNN.

98% accuracy predicting lock loss at $t = 0$, 90% at $t = -30$ s.

3-channel subset (SRCL, MICH, REFL) captures most events.

Same representation trick — wildly different physics.



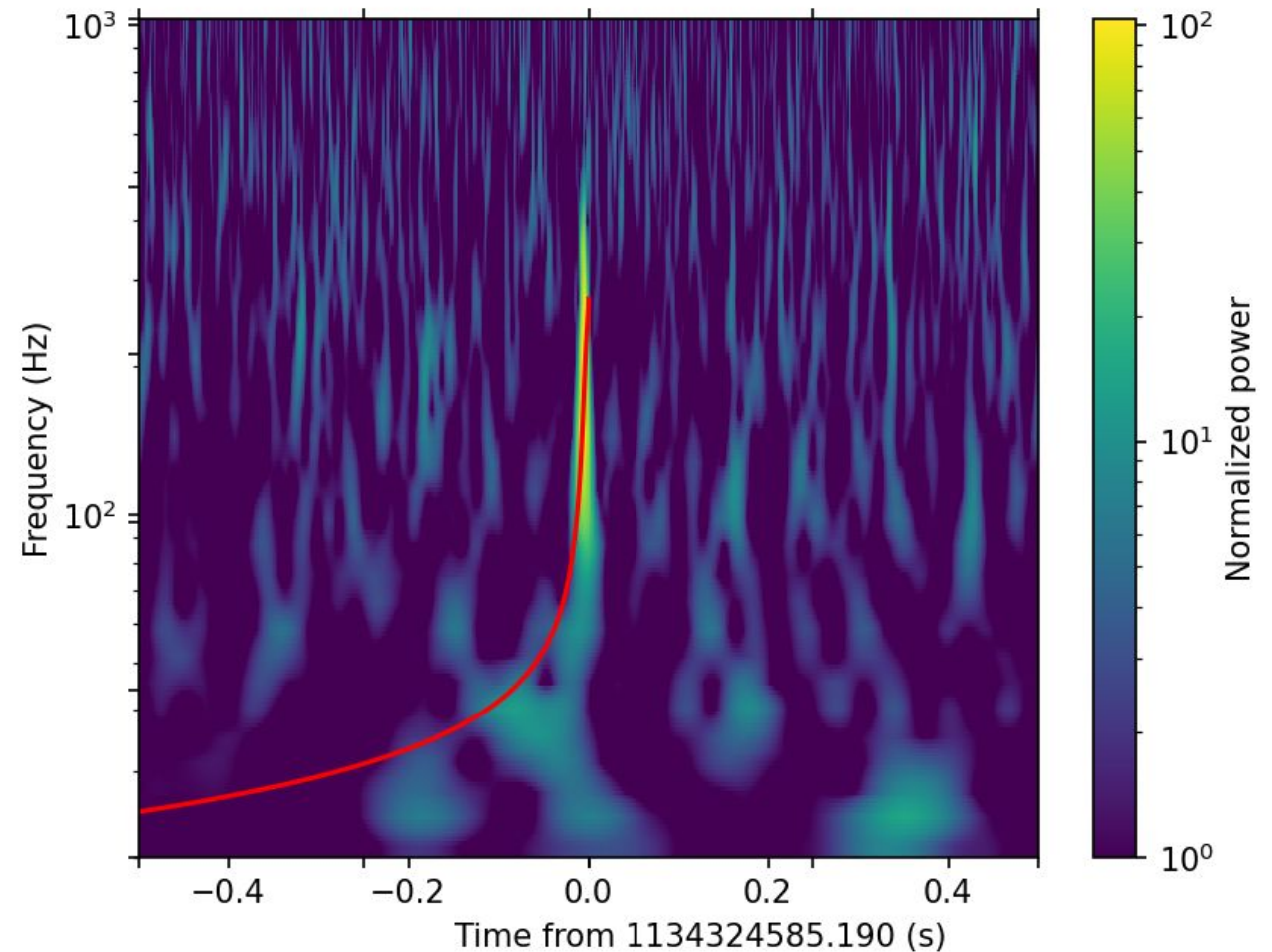
GWSkyNet — first GW real/bogus (Cabero, Mahabal, McIver 2020)

Inputs: BAYESTAR sky-map + 3D volume + detector metadata.

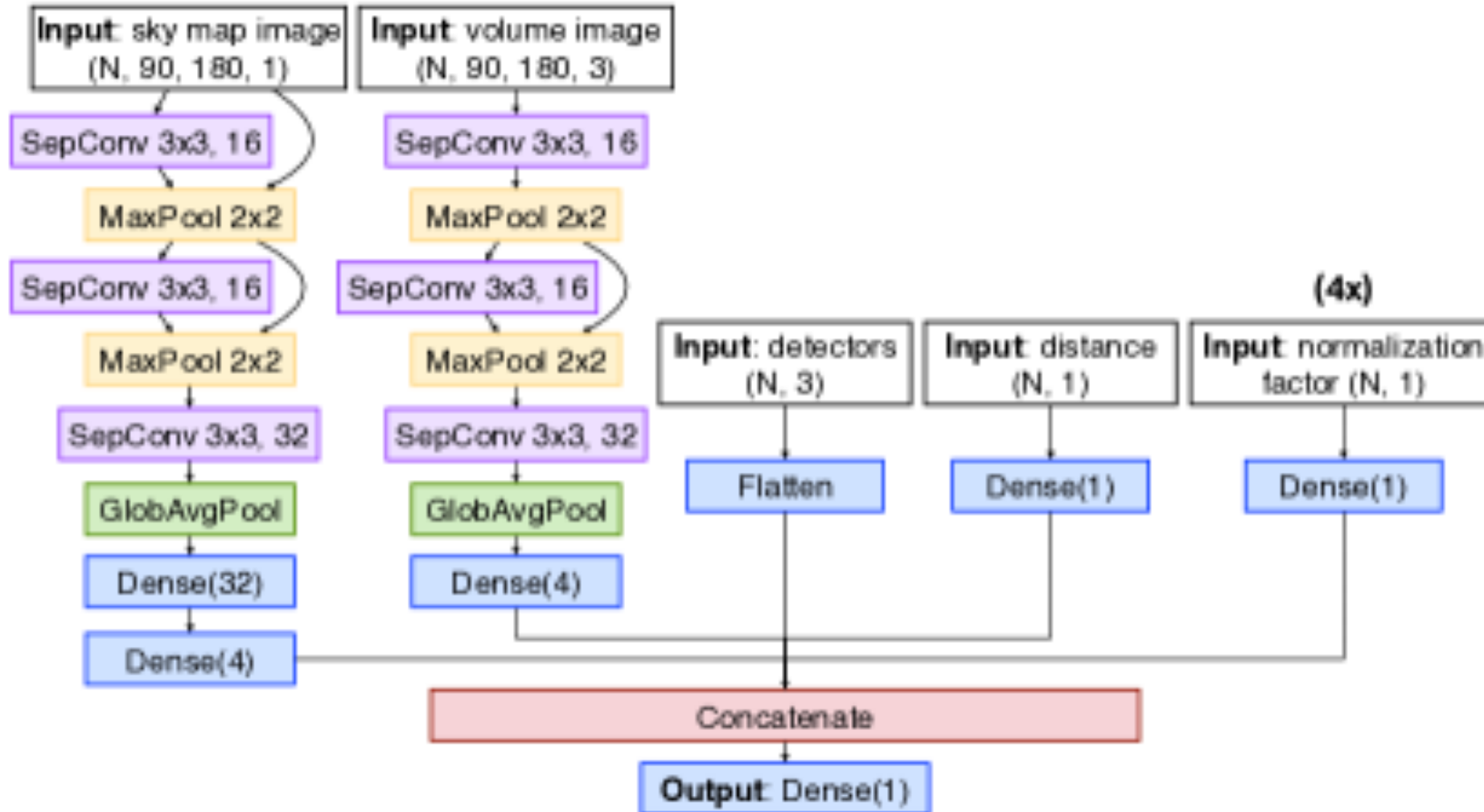
Non-sequential CNN: residual image branches + MLP for numerics.

Trained on ~2.9k simulated CBCs + 1.3k real glitches.

93.5% accuracy; 97.7% precision @ 92.4% recall.



GWSkyNet — multi-branch architecture



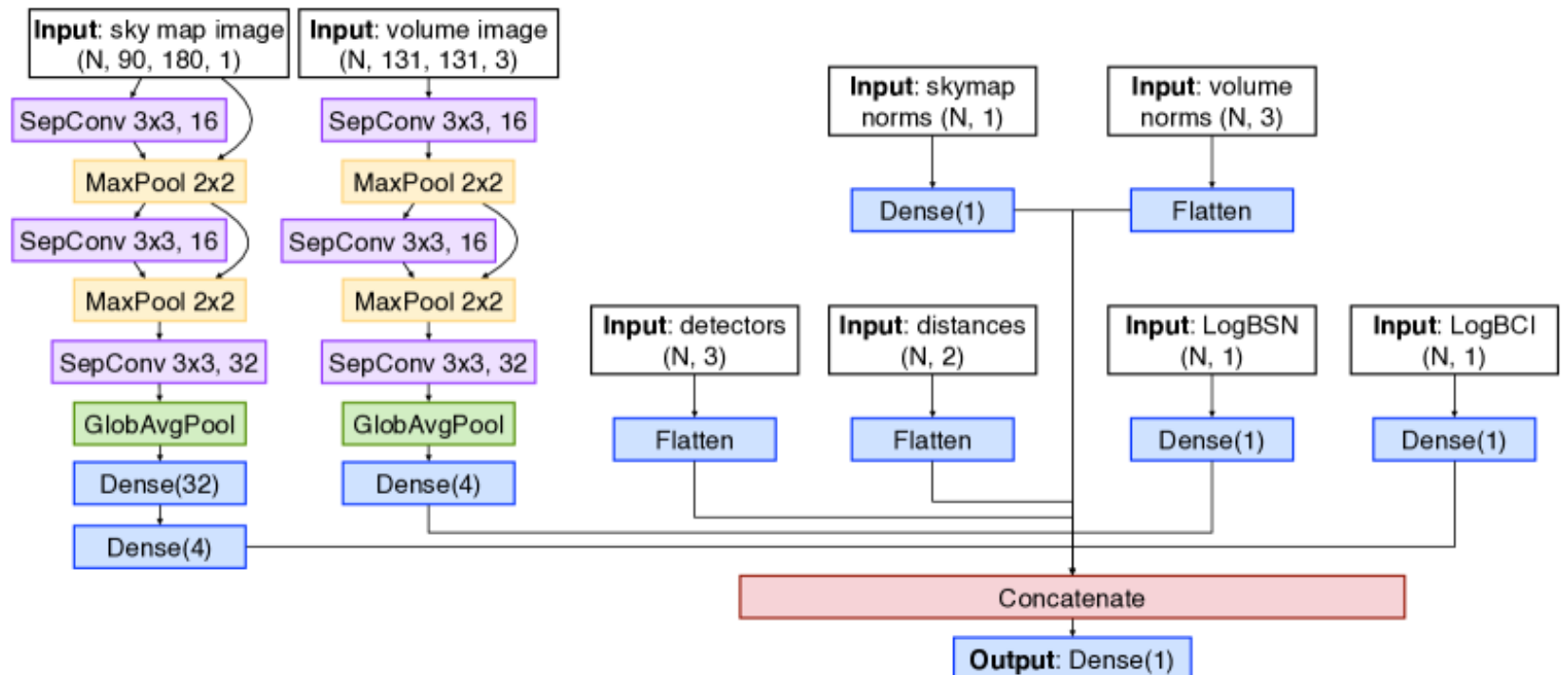
GWSkyNet-Multi — three classes (Abbott+ 2021)

LVC paper: BBH vs. NS-bearing vs. glitch.

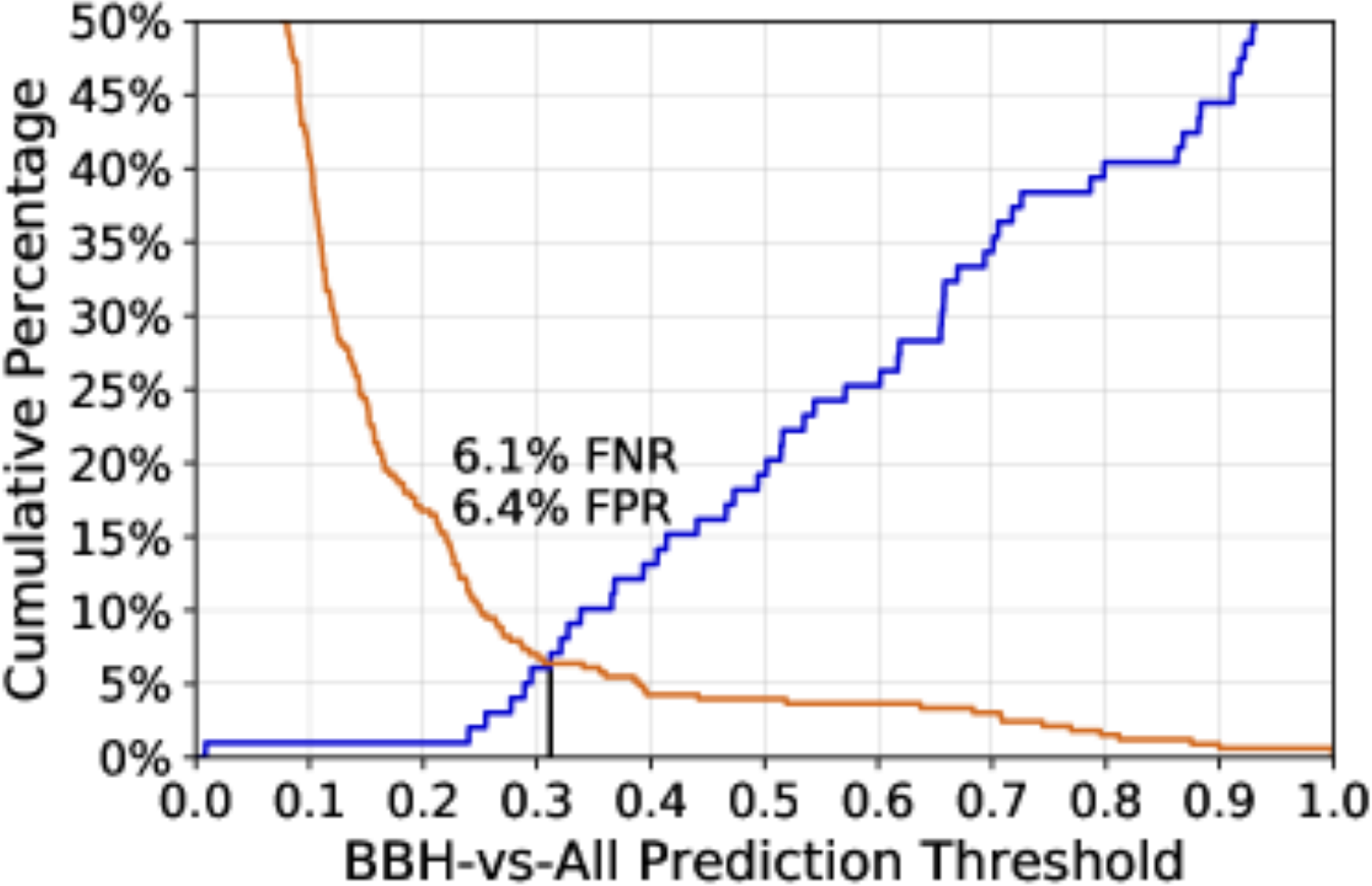
Three one-vs-all CNNs in a hierarchical decision tree.

Same inputs as GWSkyNet, broader output space.

Correctly classified 36 / 40 O3a alerts.



GWSkyNet-Multi — threshold tuning



Explainability — what is GWSkyNet looking at? (Raza+ 2023)

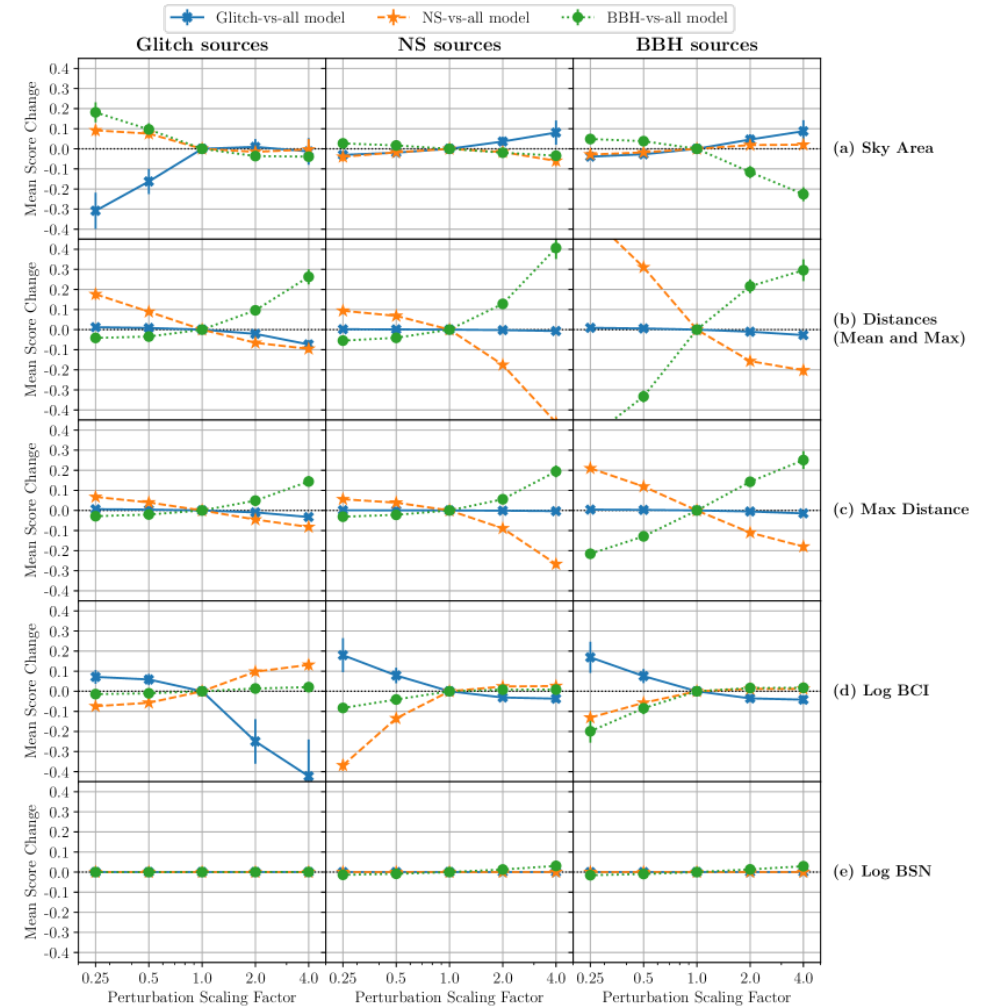
Systematic input perturbation ($0.25\times - 4\times$) to probe sensitivity.

Finding: sky-area + coherence Bayes factor drive real-vs-glitch.

Distance separates BH from NS classes.

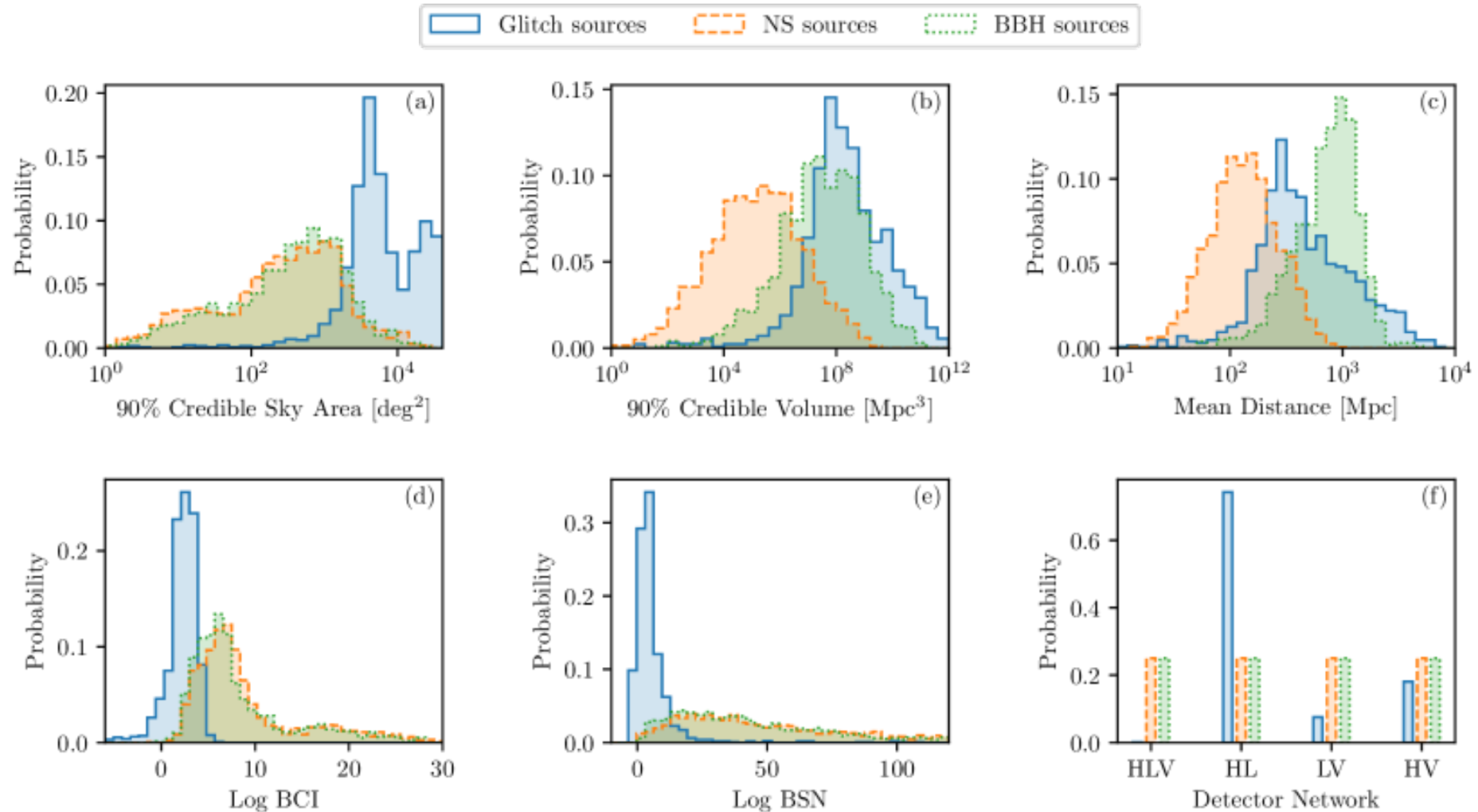
Explains 15 misclassifications in O3.

20 retrained models \rightarrow genuine error bars.



Perturbation response curves — arXiv:2308.12357

Explainability — feature distributions by class

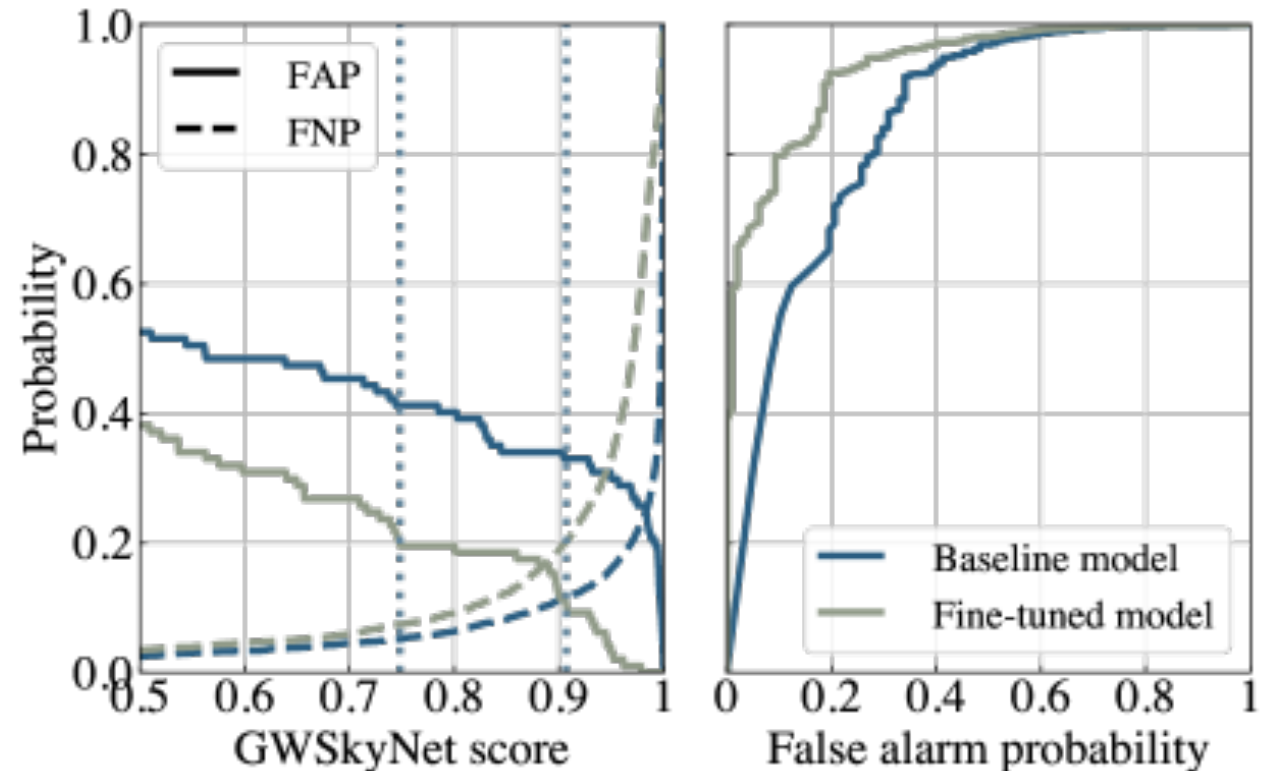


GWSkyNet II — operational in O4 (Chan+ 2024)

Retrained on O3 glitches; fine-tuned on mock-data alerts.

Transfer learning: freeze conv layers, retrain dense head.

From research code to LVK production.



GWSkyNet-Multi II — radical simplification (Raza+ 2025)

Drop image inputs; use 9 numeric summary statistics.

188 trainable parameters — 60× smaller than the predecessor.

Four classes now: glitch, BBH, NSBH, BNS.

Ensemble of 20 models → uncertainty per alert.

93% accuracy on 180 O4 alerts (May 2023 – Dec 2024).

O4 Significant Public Alerts

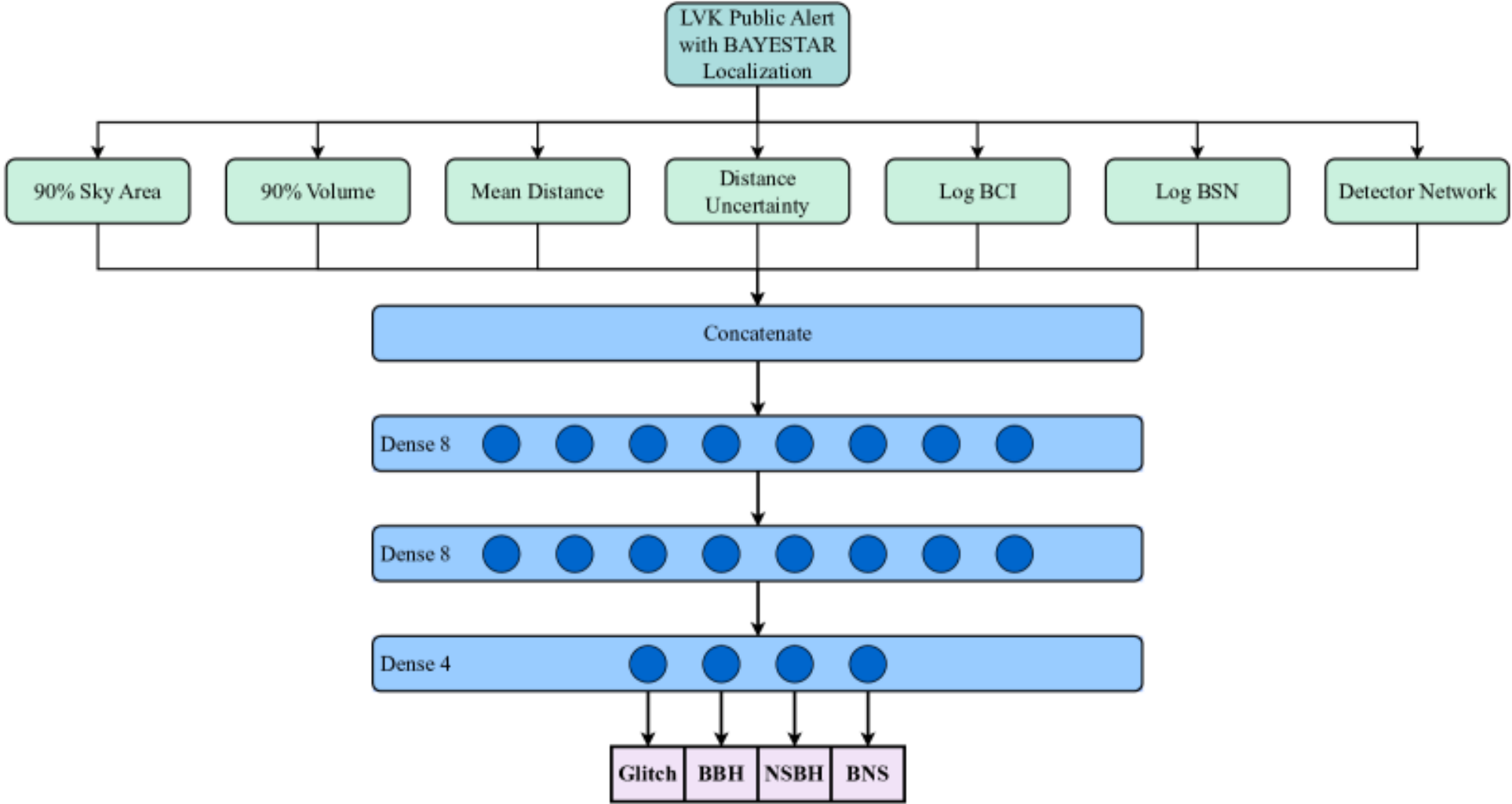
Glitch	12 (0.67)	3 (0.17)	2 (0.11)	1 (0.06)
BBH	1 (0.01)	154 (0.97)	3 (0.02)	1 (0.01)
NSBH	0 (0.00)	0 (0.00)	2 (0.67)	1 (0.33)
BNS	0 (nan)	0 (nan)	0 (nan)	0 (nan)

LVK Updated Class

Glitch BBH NSBH BNS

GWSkyNet-Multi II Predicted Class

GWSkyNet-Multi II — model architecture



9 numeric inputs → 2 dense layers → 4 class probabilities (188 params) — arXiv:2502.00297

GWSkyNet-MassGap (Raza+ 2026)

Targets the $2\text{--}5 M_{\odot}$ "lower mass gap" between heaviest NS and lightest BH.

Two outputs: $P(\text{component in mass gap})$, $P(\text{an NS is involved})$.

Best performance for chirp mass $\geq 15 M_{\odot}$; harder at low mass (needs ratio).

Mean prediction error 9% (mass-gap), 6% (NS) on O4a candidates.

arXiv:2605.00391.

GW takeaways

- One toolbox, a growing family of tasks.
- From "is it real?" → "what is it?" → "where on the mass spectrum?"
- Real-time alerts + offline inference both matter.
- DMDT from variability shows up here too.
- Models get *simpler* as understanding improves (GWSkyNet-Multi II: 188 params).

5. Radio transients

FRBs and the candidate-rate problem

CHIME / DSA-110 generate thousands of dispersed candidates per day.

Dynamic spectra (time \times frequency) are 2-D — natural CNN inputs.

More than real/bogus: morphology drives follow-up choices.

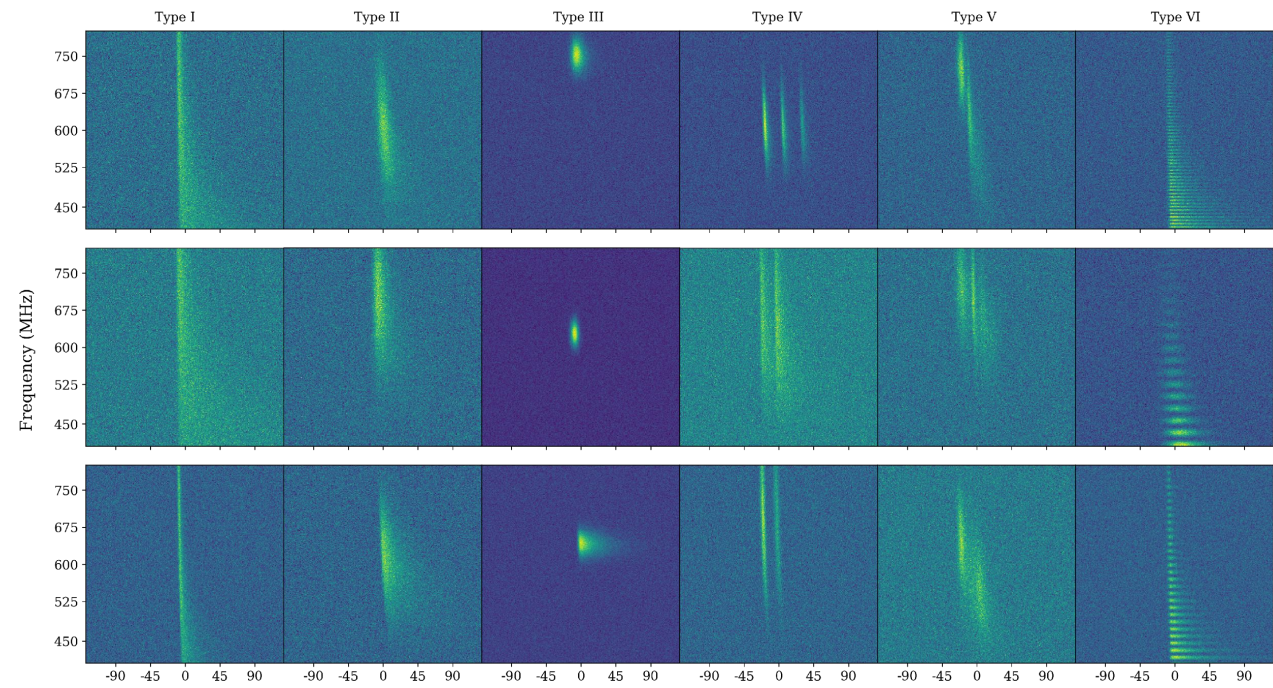
Repeaters look different from one-offs — and we want to know which is which fast.

Frabjous — DL for FRB morphology (Kumar+ 2025)

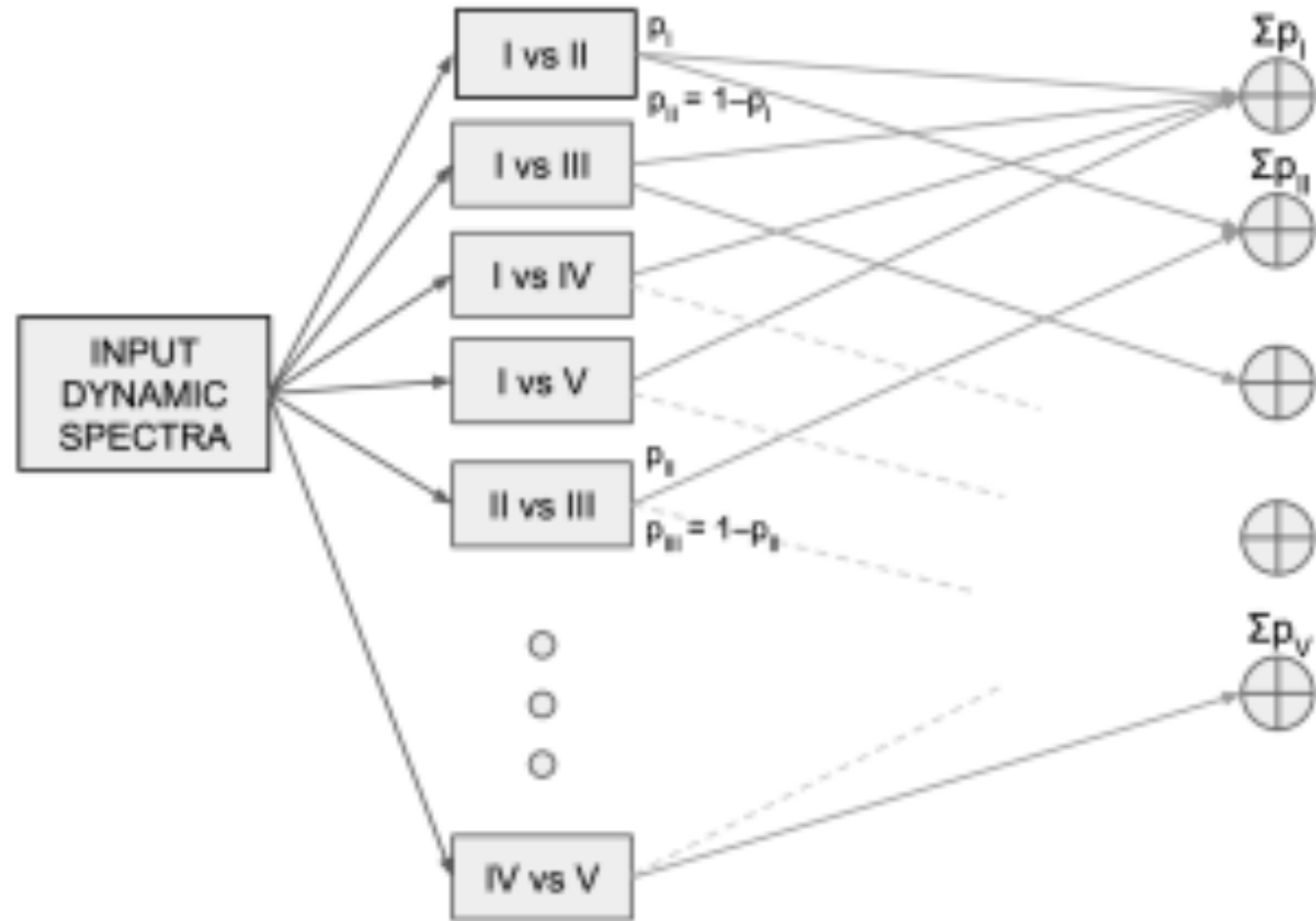
Six morphology archetypes (broadband, narrowband, multi-component, drifting, ...).

10 binary classifiers \oplus confidence-voting matrix \rightarrow 5-way label.

55% accuracy on real CHIME catalog (random baseline: 20%).



Frabjous — classifier ensemble + voting matrix



6. Language models

Why LLMs for astronomy?

Literature grows faster than humans can read.

Codebases, archives, manuals — all natural-language interfaces.

Multi-modal: text + light curves + images + spectra.

Agents that vet alerts, plan follow-up, draft reports.

But — can we trust them? Or, how much can we trust them?

AstroAlertBench (Chen+ 2026)

Benchmark of multimodal LLMs on real ZTF alerts (arXiv:2605.05573).

1,500 alerts × 13 frontier models. (1,500 is tiny)

Three evaluation stages:

metadata grounding → scientific reasoning → hierarchical classification.

High accuracy does ****not**** track "honesty" (self-evaluation of reasoning).

1. First-alert inputs

Image Triplet + Alert Metadata



Science

Reference

Difference

Representative Metadata

objectId	ZTF20example
fid	2 (r band)
magpsf	18.45
sgscore1	0.99

Three 63×63 cutouts centered on candidate position, plus prompt-facing fields extracted from the alert packet.

Serialize

2. Prompt Construction

A Consistent Multimodal Prompt

Alert ID: ZTF20example

Images show a localized candidate residual.

Metadata: r band;
magpsf = 18.45;
sgscore1 = 0.99.

Task: Answer Parts A-C
in the required format.

Includes:

- Image triplets
- Metadata
- Field definitions
- Output instructions

Evaluate

3. Structured Model Response

Machine-parseable Parts A-C

VLM

Part A

filter_band: r
subtraction_sign: positive

Part B

leading_interpretation:
supernova candidate

Part C

stage1: real_object
stage2: astrophysical
stage3: supernova

Scored on metadata
grounding, rationale, and
staged classification.

Cross-cutting lessons

Representation > architecture: DMDT, Q-transforms, image triplets.

Labels — not compute — are the real bottleneck. Active learning is everywhere.

Asymmetric error costs shape the loss (missed kilonova vs. extra alert).

Humans-in-the-loop are everywhere.

What's next — the LSST / multi-messenger era

LSST: $\sim 10\times$ ZTF rates for a decade, starting soon.

Foundation models across surveys and wavebands.

Joint EM + GW + neutrino inference, not separate pipelines.

Real-time decisions delegated to (audited) agents.

The toolbox built on ZTF + LIGO + TESS + CHIME generalises.

Interpretability and explainability

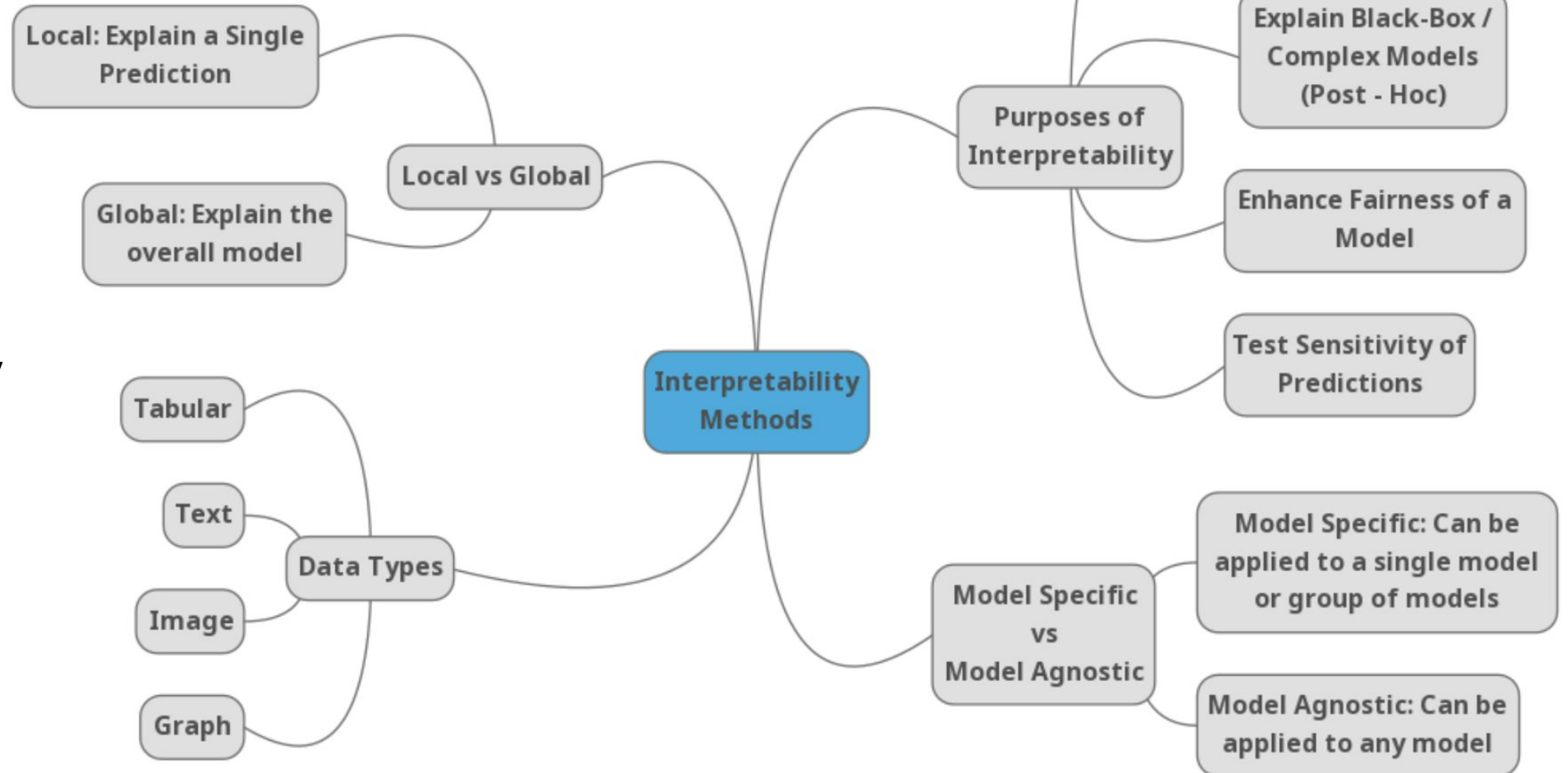
Linardatos,
Papastefanopoulos,
Kotsiantis 2020

Post-hoc

LIME: Local
Interpretable
Model-agnostic
Explanations

SHAP: Shapley
Additive
Explanations

Fairness



References

Images: 1907.11259, 1904.05920, 1710.01422

Light curves: 1709.06257, 2009.14071, 2102.11304, 2312.00143,
2101.09227

Spectra: 2104.12980, 2412.08601

GW: 1910.12143, 2010.11829, 2111.04015, 2308.12357, 2408.06491,
2502.00297, 2605.00391

Radio: 2507.14854

LLMs: 2605.05573