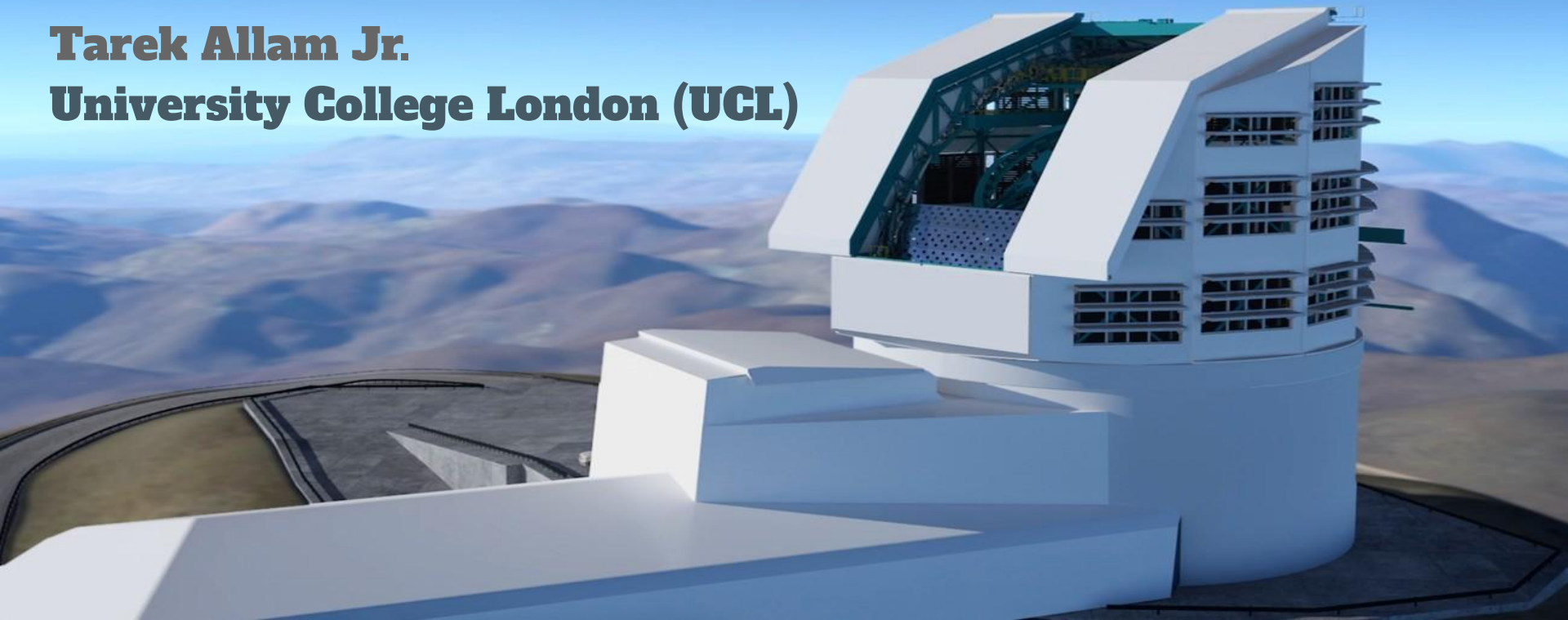# Low-latency High-throughput Classification using Deep Model Compression
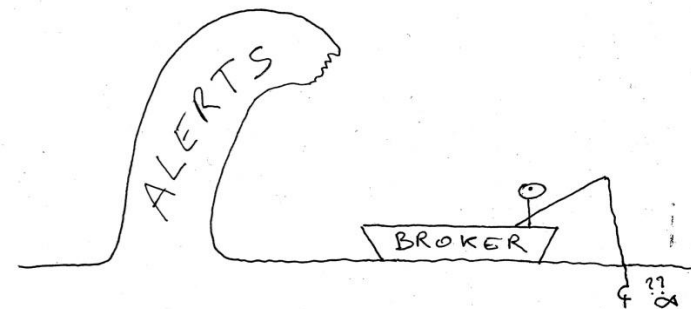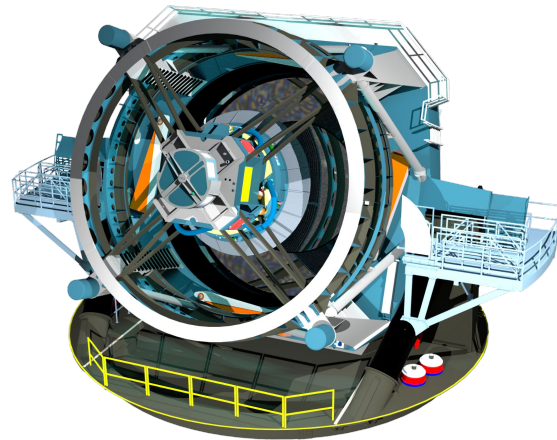
**Tarek Allam Jr.**
**University College London (UCL)**

# Motivation: A Deluge of Data

**Overview**

- 10 million alerts, per night!

- Machine Learning methods are now critical

- Accurate and fast classification required for follow up

- Desire for fast re-training of models



Peloton et al. 2020

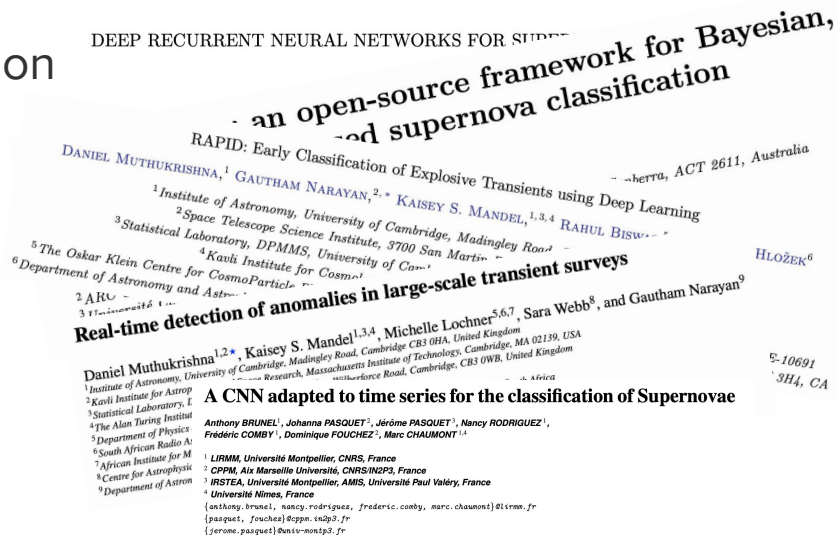# Deep Learning to the Rescue!?

**The death of feature engineering?**

- Exploiting inherent time-series information

**RNNs** (inc. SuperNNova, RAPID)

**CNNs** (inc. TCN)

**Transformers** (inc. t2)

● ● ●

DEEP RECURRENT NEURAL NETWORKS FOR SUPE...

an open-source framework for Bayesian,
...ed supernova classification

RAPID: Early Classification of Explosive Transients using Deep Learning

DANIEL MUTHUKRISHNA,[1] GAUTHAM NARAYAN,[2,*] KAISEY S. MANDEL,[1,3,4] RAHUL BIS...

[1] Institute of Astronomy, University of Cambridge, Madingley Road
[2] Space Telescope Science Institute, 3700 San Martin...
[3] Statistical Laboratory, DPMMS, University of Cam...
[4] Kavli Institute for Cosmo...
[5] The Oskar Klein Centre for CosmoParticle...
[6] Department of Astronomy and Astro...

...erra, ACT 2611, Australia

HLOŽEK[6]

Real-time detection of anomalies in large-scale transient surveys

Daniel Muthukrishna[1,2], Kaisey S. Mandel[1,3,4], Michelle Lochner[5,6,7], Sara Webb[8], and Gautham Narayan[9]

[1] Institute of Astronomy, University of Cambridge, Madingley Road, Cambridge CB3 0HA, United Kingdom
[2] Kavli Institute for Astrop... Massachusetts Institute of Technology, Cambridge, MA 02139, USA
[3] Statistical Laboratory, ...
[4] The Alan Turing Institut...
[5] Department of Physics ...berforce Road, Cambridge, CB3 0WB, United Kingdom
[6] South African Radio A...
[7] African Institute for M...
[8] Centre for Astrophysic...
[9] Department of Astron...

...a Africa

...5-10691
...3H4, CA

A CNN adapted to time series for the classification of Supernovae

Anthony BRUNEL[1], Johanna PASQUET[2], Jérôme PASQUET[3], Nancy RODRIGUEZ[1], Frédéric COMBY[1], Dominique FOUCHEZ[2], Marc CHAUMONT[1,4]

[1] LIRMM, Université Montpellier, CNRS, France
[2] CPPM, Aix Marseille Université, CNRS/IN2P3, France
[3] IRSTEA, Université Montpellier, AMIS, Université Paul Valéry, France
[4] Université Nîmes, France

{anthony.brunel, nancy.rodriguez, frederic.comby, marc.chaumont}@lirmm.fr
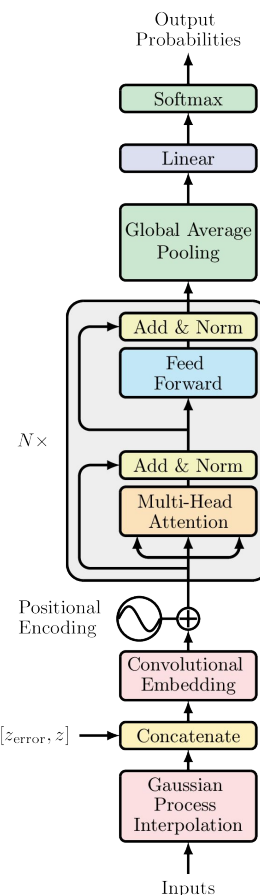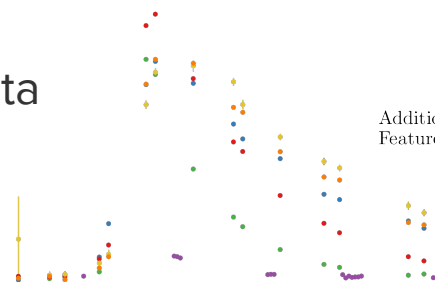{pasquet, fouchez}@cppm.in2p3.fr
{jerome.pasquet}@univ-montp3.fr

# The Time-Series Transformer [t2]
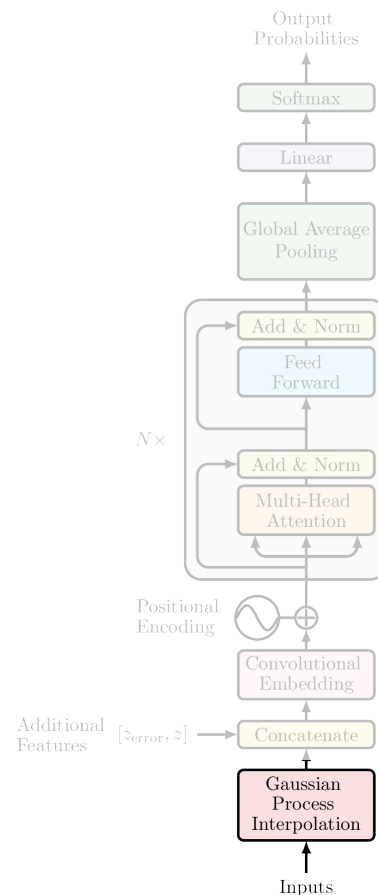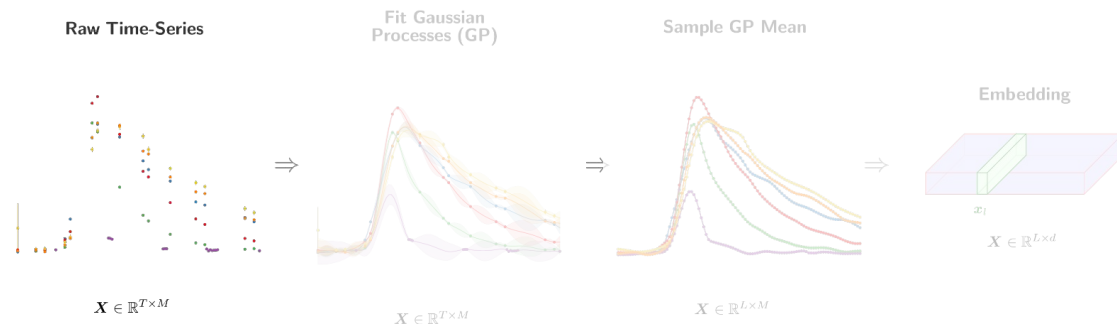
**Encoder++**

- *Global Average Pooling* to allow for Class Activation Maps

- *Convolutional Embedding* maps time-series into a vector space

- *Concatenate* additional features

- *GP Interpolation* to handle irregular data
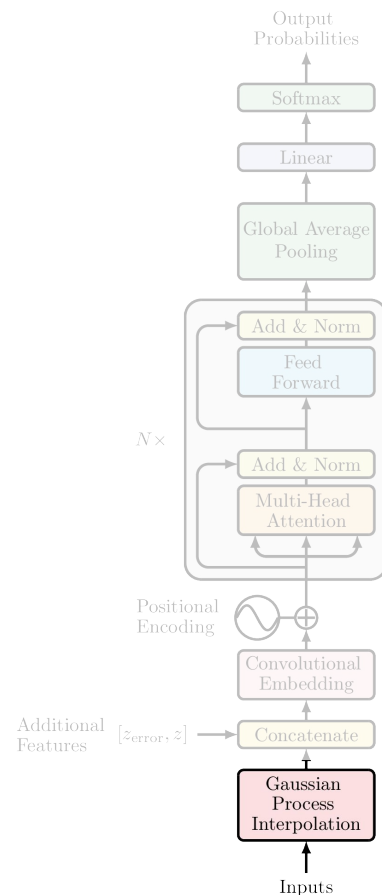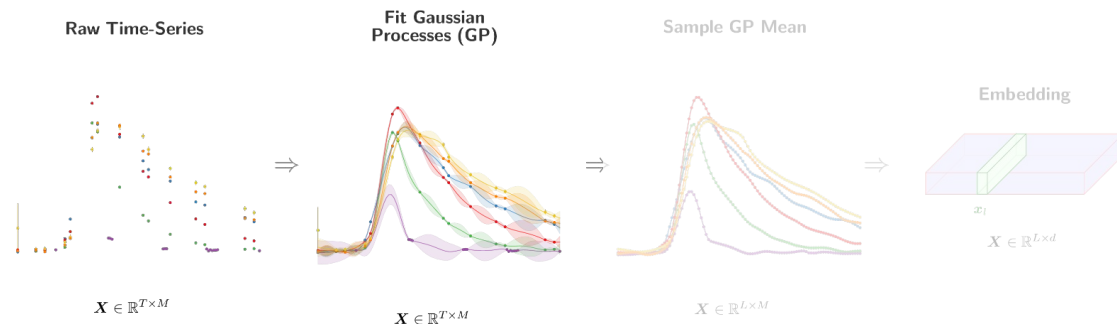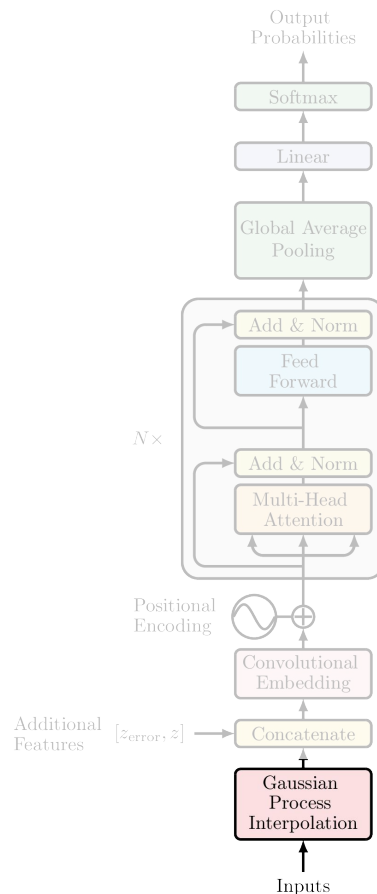
# ML Pipeline

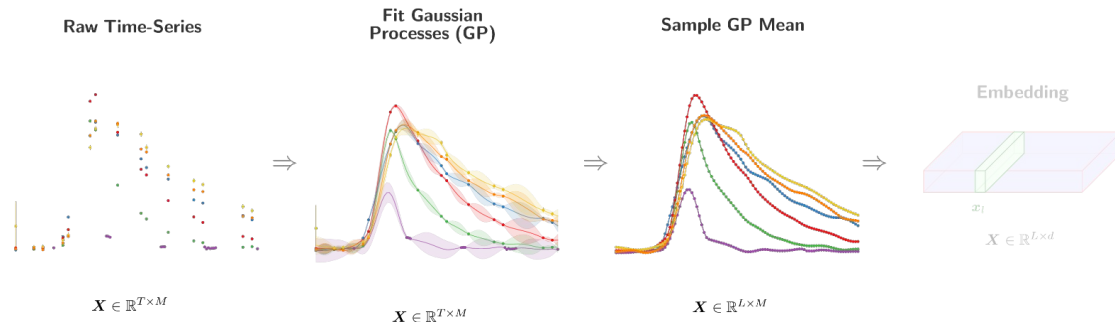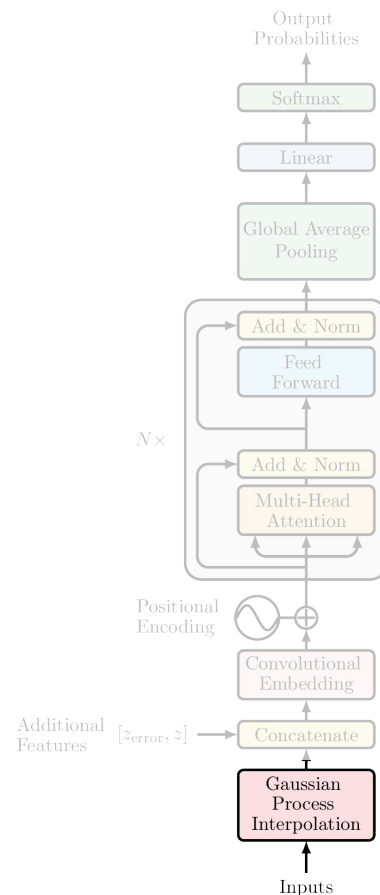## Gaussian Process Interpolation

- 2D-Matern kernel for Gaussian Process interpolation
- Evaluate at regular period, 100 points in our case

# ML Pipeline

## Gaussian Process Interpolation

- 2D-Matern kernel for Gaussian Process interpolation
- Evaluate at regular period, 100 points in our case



**Raw Time-Series**

$$\boldsymbol{X} \in \mathbb{R}^{T \times M}$$

**Fit Gaussian Processes (GP)**

$$\boldsymbol{X} \in \mathbb{R}^{T \times M}$$

**Sample GP Mean**

$$\boldsymbol{X} \in \mathbb{R}^{L \times M}$$

**Embedding**

$$x_l$$

$$\boldsymbol{X} \in \mathbb{R}^{L \times d}$$

Output Probabilities

Softmax

Linear

Global Average Pooling

Add & Norm

Feed Forward

$N \times$

Add & Norm

Multi-Head Attention

Positional Encoding

Convolutional Embedding

Additional Features $[z_{error}, z]$ → Concatenate

Gaussian Process Interpolation

Inputs

# ML Pipeline

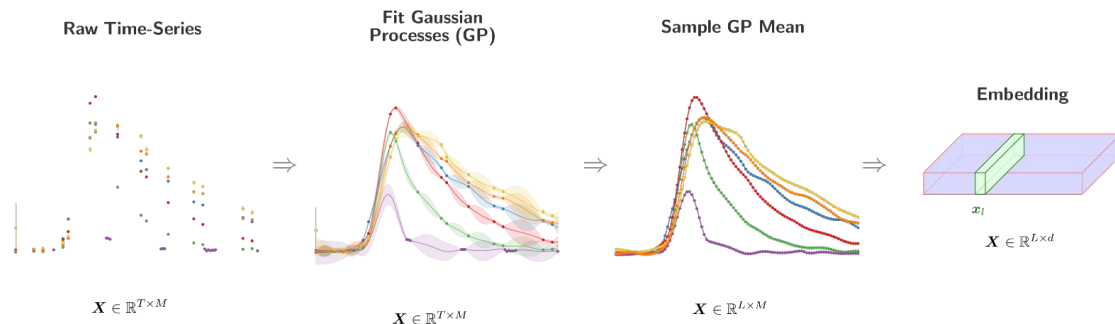## Gaussian Process Interpolation

- 2D-Matern kernel for Gaussian Process interpolation
- Evaluate at regular period, 100 points in our case

# ML Pipeline

## Gaussian Process Interpolation

- 2D-Matern kernel for Gaussian Process interpolation
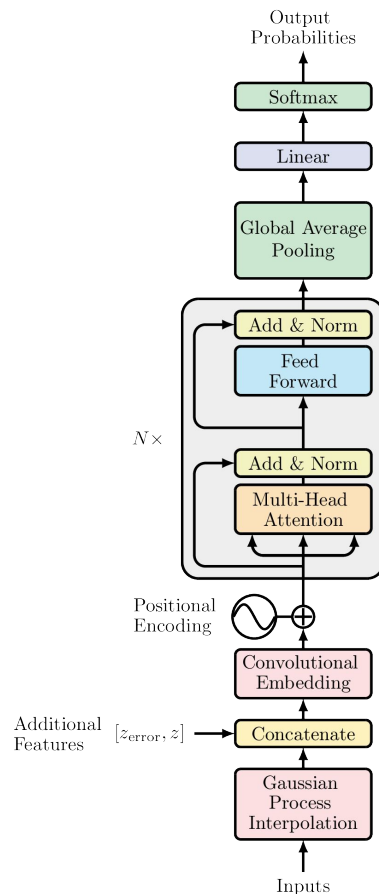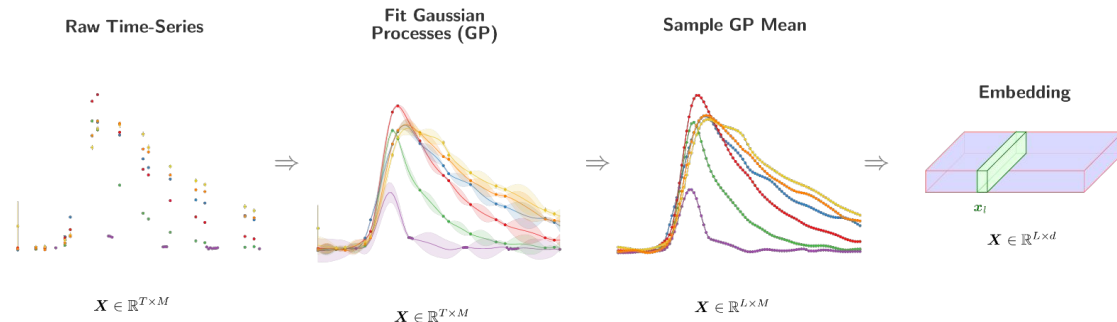- Evaluate at regular period, 100 points in our case

# ML Pipeline

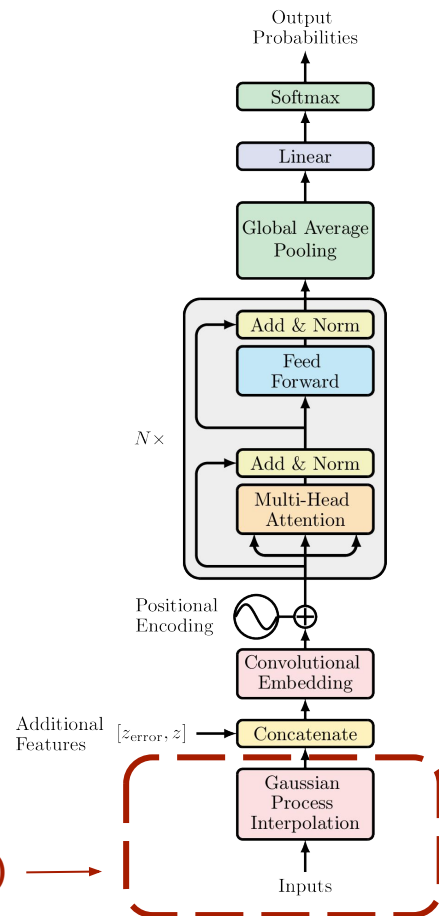## Gaussian Process Interpolation

- 2D-Matern kernel for Gaussian Process interpolation
- Evaluate at regular period, 100 points in our case

## Deep Learning Magic ...



**Raw Time-Series**

$\boldsymbol{X} \in \mathbb{R}^{T \times M}$

**Fit Gaussian Processes (GP)**

$\boldsymbol{X} \in \mathbb{R}^{T \times M}$

**Sample GP Mean**

$\boldsymbol{X} \in \mathbb{R}^{L \times M}$

**Embedding**

$\boldsymbol{x}_l$

$\boldsymbol{X} \in \mathbb{R}^{L \times d}$



Output Probabilities

Softmax

Linear

Global Average Pooling

Add & Norm

Feed Forward

$N\times$

Add & Norm

Multi-Head Attention

Positional Encoding

Convolutional Embedding

Additional Features $[z_{\text{error}}, z]$ → Concatenate
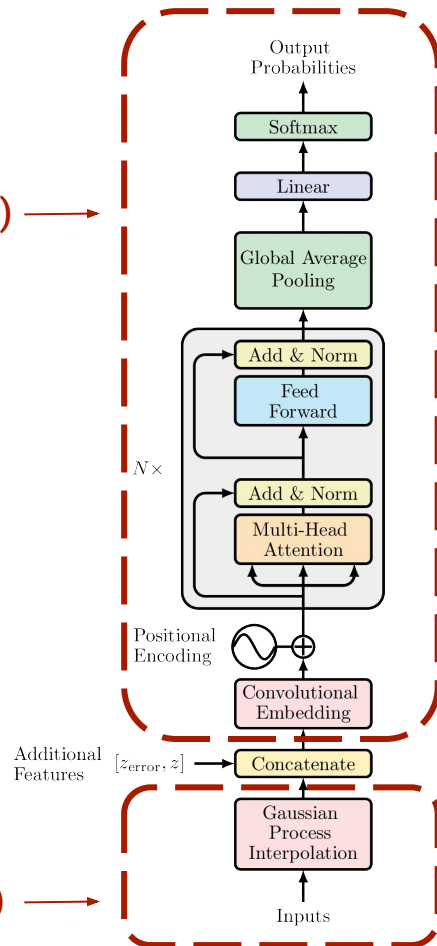
Gaussian Process Interpolation

Inputs

# Understanding Bottlenecks

# Understanding Bottlenecks



`model.predict(processed_alert)` →
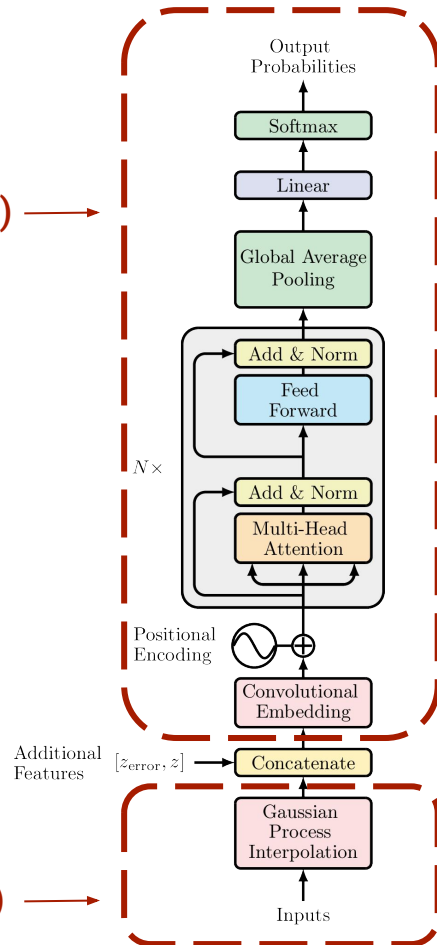
`fit_gps(raw_alert)` →

# Understanding Bottlenecks

```
Total time: 5.85664 s
File: get_models.py
Function: get_model at line 29

Total time: 1.47076 s
File: ztf-load-run-lpa.py
Function: t2_probs at line 55

Line #    Hits        Time Per Hit \% Time Line Contents
==============================================================
...
206       16      139.6    8.7    9.5     df_gp_mean = generate_gp_all_objects()
...
...
...
...
...
212       8       180.8   22.6   12.3     X = df_gp_mean[cols]
213       8        12.3    1.5    0.8     X = rs(X)
...
...
...
217       8      1101.7  137.7   74.9     y_preds = model.predict(X)
```

`model.predict(processed_alert)` →

`fit_gps(raw_alert)` →

# Understanding Bottlenecks

👇

```
Total time: 5.85664 s
File: get_models.py
Function: get_model at line 29


Total time: 1.47076 s
File: ztf-load-run-lpa.py
Function: t2_probs at line 55


Line #    Hits        Time Per Hit \% Time Line Contents
==============================================================
...
206       16     139.6    8.7    9.5    df_gp_mean = generate_gp_all_objects()
...
...
...
...
...
212       8      180.8   22.6   12.3    X = df_gp_mean[cols]
213       8       12.3    1.5    0.8    X = rs(X)
...
...
...
217       8     1101.7  137.7   74.9    y_preds = model.predict(X)
```
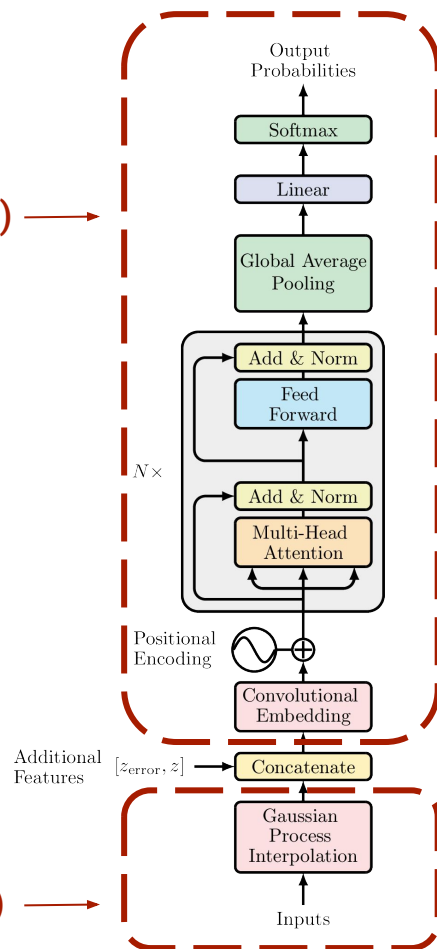
model.predict(processed_alert) →

fit_gps(raw_alert) →

# Understanding Bottlenecks

```
Total time: 5.85664 s
File: get_models.py
Function: get_model at line 29


Total time: 1.47076 s
File: ztf-load-run-lpa.py
Function: t2_probs at line 55


Line #    Hits      Time Per Hit \% Time Line Contents
================================================================
...
206       16        139.6    8.7    9.5    df_gp_mean = generate_gp_all_objects()
...
...
...
...
...
212       8         180.8    22.6   12.3   X = df_gp_mean[cols]
213       8         12.3     1.5    0.8    X = rs(X)
...
...
...
217       8         1101.7   137.7  74.9   y_preds = model.predict(X)
```
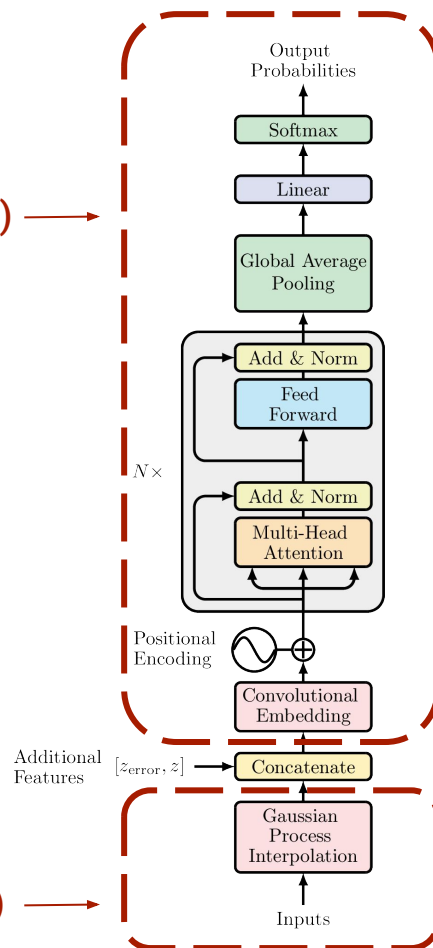
model.predict(processed_alert) →
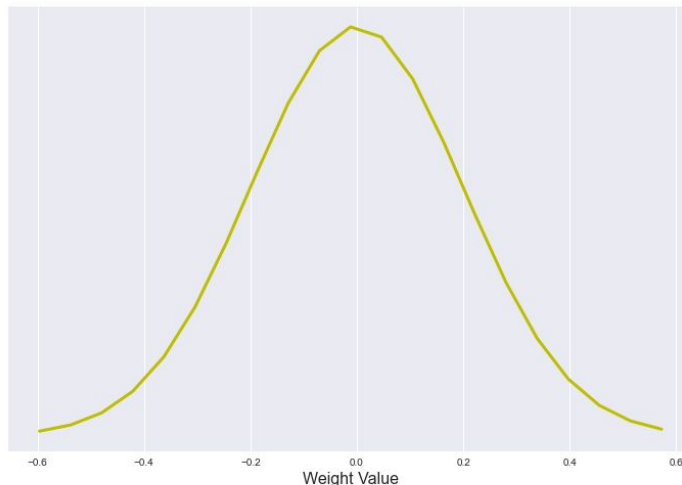
fit_gps(raw_alert) →

# Deep Model Compression: An *Almost* Free Lunch

- Standout paper: Han et al. 2015: *Deep Compression*

  - Clustering

  - Pruning

  - Quantization

- Lossy process – expect some loss in accuracy (but improved compute and storage costs!)

# Deep Model Compression: Clustering

**Weight Clustering**

- Reduce number of unique weights with shared weight values

- Can then leverage Huffman encoding

- Note: best to avoid clustering early layers!

# Deep Model Compression: Clustering

**Weight Clustering**

- Reduce number of unique weights with shared weight values

- Can then leverage Huffman encoding

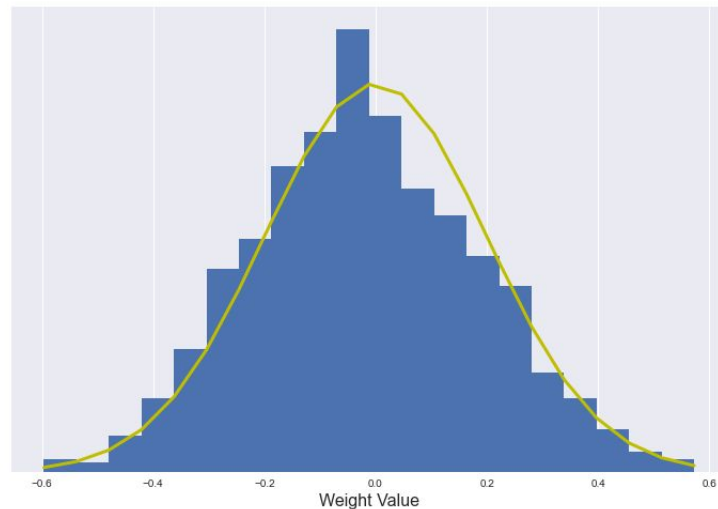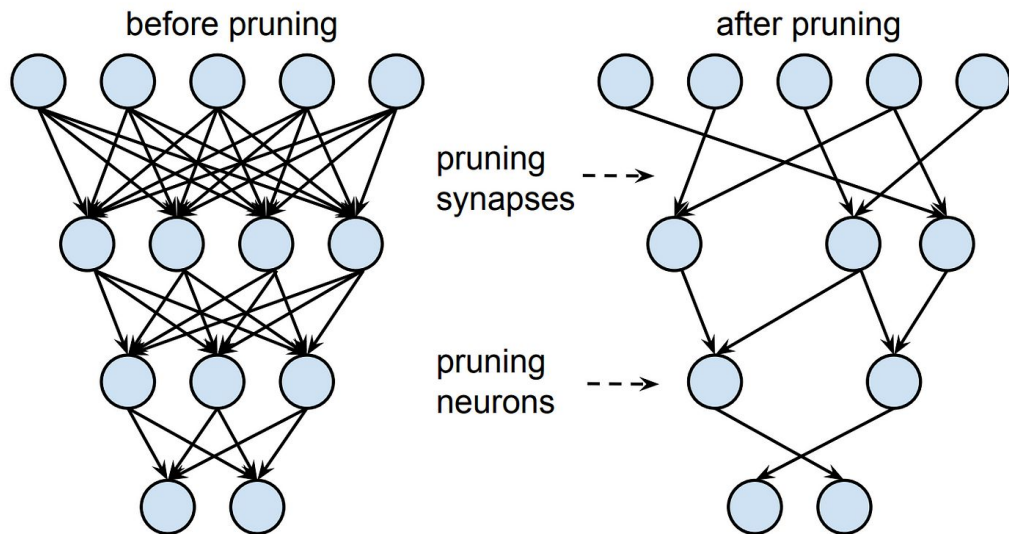- Note: best to avoid clustering early layers!

# Deep Model Compression: Pruning

**Weight Pruning**
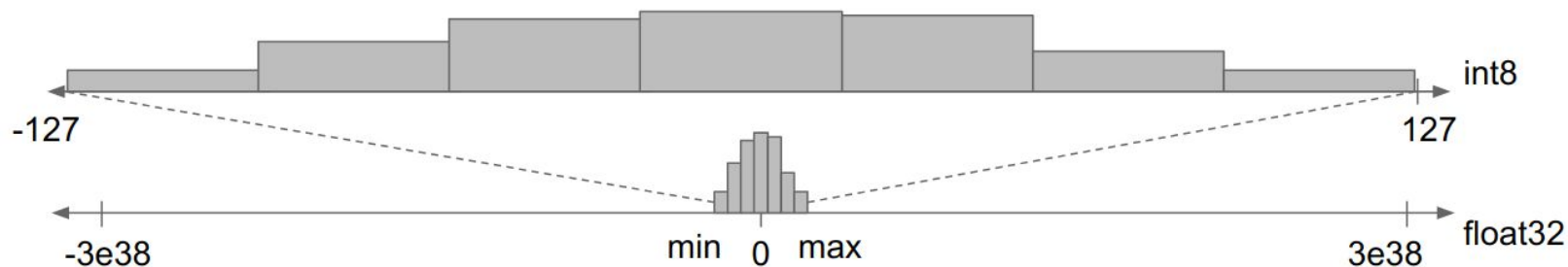
- Set low magnitude weights to zero

- Can then leverage Huffman encoding (again)

- Done *during* training

# Deep Model Compression: Quantization

**Weight Quantization**

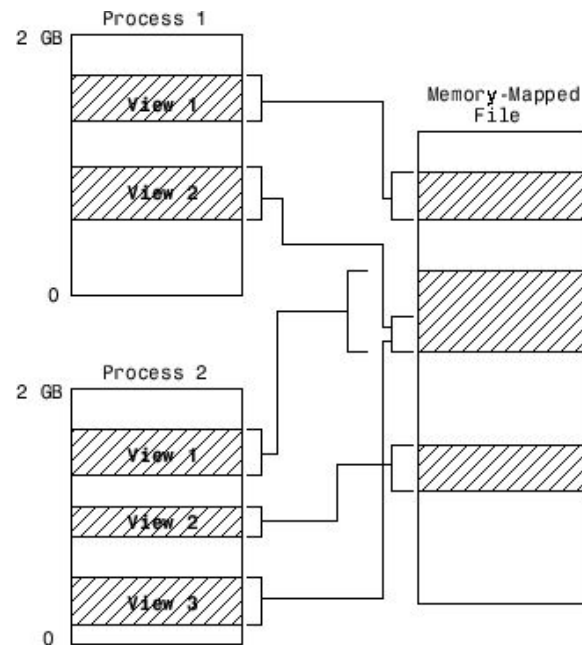- Reduce precision of stored weight values

- Computations are still done at `float32`

# Efficient Frameworks & File Formats

**From TensorFlow to TensorFlow-Lite**
- Operator Fusion
- `ProtocolBuffers → FlatBuffers`

# Local Tests

| Compression Method | Model Size (kb) | Load Latency ($s^{-3}$) | Inference Latency (s) | Loss |
|---|---|---|---|---|
| Baseline | 1100 | 6324.145 | 0.333 | 0.968 |

# Local Tests

| Compression Method | Model Size (kb) | Load Latency ($s^{-3}$) | Inference Latency (s) | Loss |
|---|---|---|---|---|
| Baseline | 1100 | 6324.145 | 0.333 | 0.968 |
| Baseline + Huffman | 244 | 6015.565 | 0.224 | 0.968 |

# Local Tests

| Compression Method | Model Size (kb) | Load Latency ($s^{-3}$) | Inference Latency (s) | Loss |
|---|---|---|---|---|
| Baseline | 1100 | 6324.145 | 0.333 | 0.968 |
| Baseline + Huffman | 244 | 6015.565 | 0.224 | 0.968 |
| Clustering | 892 | 5559.868 | 0.227 | 0.836 |

# Local Tests

| Compression Method | Model Size (kb) | Load Latency ($s^{-3}$) | Inference Latency (s) | Loss |
|---|---|---|---|---|
| Baseline | 1100 | 6324.145 | 0.333 | 0.968 |
| Baseline + Huffman | 244 | 6015.565 | 0.224 | 0.968 |
| Clustering | 892 | 5559.868 | 0.227 | 0.836 |
| Clustering + Pruning | 688 | 5721.021 | 0.230 | 1.017 |

# Local Tests

| COMPRESSION METHOD | MODEL SIZE (KB) | LOAD LATENCY ($s^{-3}$) | INFERENCE LATENCY (S) | LOSS |
|---|---|---|---|---|
| BASELINE | 1100 | 6324.145 | 0.333 | 0.968 |
| BASELINE + HUFFMAN | 244 | 6015.565 | 0.224 | 0.968 |
| CLUSTERING | 892 | 5559.868 | 0.227 | 0.836 |
| CLUSTERING + PRUNING | 688 | 5721.021 | 0.230 | 1.017 |
| CLUSTERING + HUFFMAN | 240 | 4991.857 | 0.223 | 0.836 |

# Local Tests

| Compression Method | Model Size (kB) | Load Latency ($s^{-3}$) | Inference Latency (s) | Loss |
|---|---|---|---|---|
| Baseline | 1100 | 6324.145 | 0.333 | 0.968 |
| Baseline + Huffman | 244 | 6015.565 | 0.224 | 0.968 |
| Clustering | 892 | 5559.868 | 0.227 | 0.836 |
| Clustering + Pruning | 688 | 5721.021 | 0.230 | 1.017 |
| Clustering + Huffman | 240 | 4991.857 | 0.223 | 0.836 |
| Clustering + Pruning + Huffman | 128 | 5251.288 | 0.228 | 1.017 |

# Local Tests

| Compression Method | Model Size (kb) | Load Latency ($s^{-3}$) | Inference Latency (s) | Loss |
|---|---|---|---|---|
| Baseline | 1100 | 6324.145 | 0.333 | 0.968 |
| Baseline + Huffman | 244 | 6015.565 | 0.224 | 0.968 |
| Clustering | 892 | 5559.868 | 0.227 | 0.836 |
| Clustering + Pruning | 688 | 5721.021 | 0.230 | 1.017 |
| Clustering + Huffman | 240 | 4991.857 | 0.223 | 0.836 |
| Clustering + Pruning + Huffman | 128 | 5251.288 | 0.228 | 1.017 |
| †Clustering | 92 | 0.426 | 0.046 | 0.836 |

# Local Tests

| Compression Method | Model Size (kB) | Load Latency ($s^{-3}$) | Inference Latency (s) | Loss |
|---|---|---|---|---|
| Baseline | 1100 | 6324.145 | 0.333 | 0.968 |
| Baseline + Huffman | 244 | 6015.565 | 0.224 | 0.968 |
| Clustering | 892 | 5559.868 | 0.227 | 0.836 |
| Clustering + Pruning | 688 | 5721.021 | 0.230 | 1.017 |
| Clustering + Huffman | 240 | 4991.857 | 0.223 | 0.836 |
| Clustering + Pruning + Huffman | 128 | 5251.288 | 0.228 | 1.017 |
| †Clustering | 92 | 0.426 | 0.046 | 0.836 |
| **†Clustering + Quantization** | **60** | **0.271** | **0.043** | **0.834** |

# Local Tests

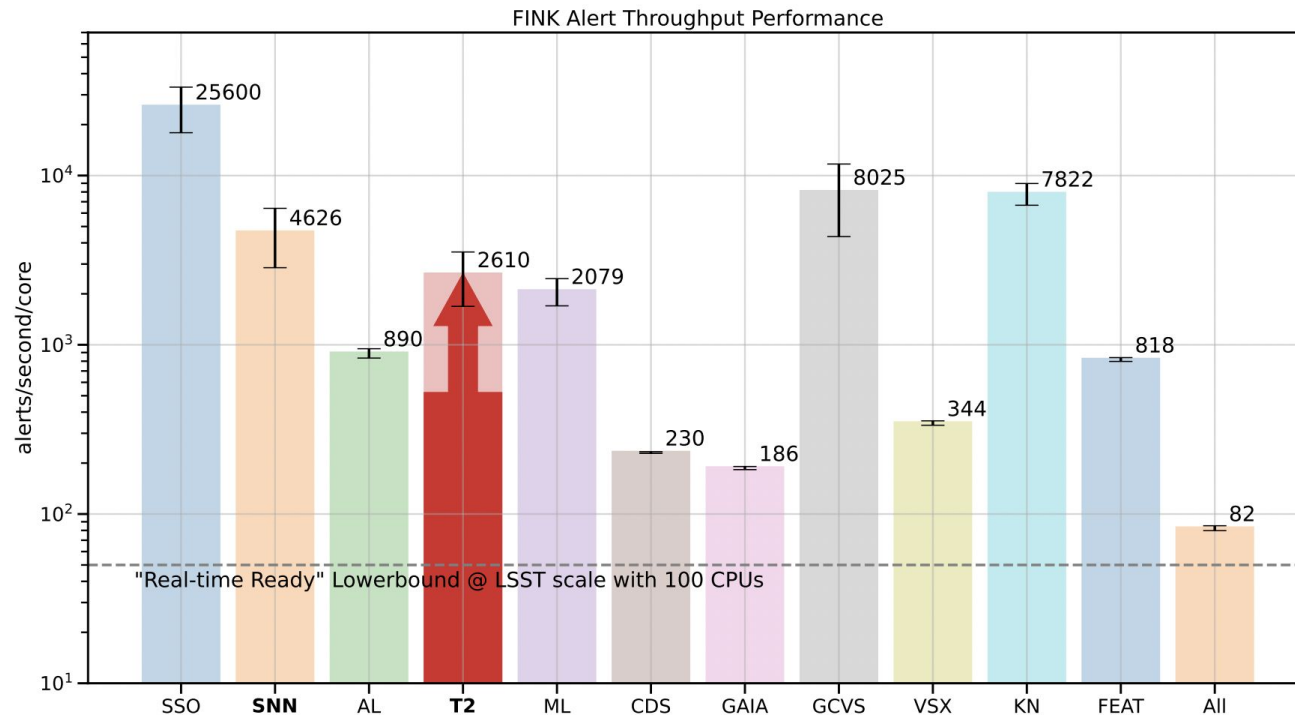| Compression Method | Model Size (kB) | Load Latency ($s^{-3}$) | Inference Latency (s) | Loss |
|---|---|---|---|---|
| Baseline | 1100 | 6324.145 | 0.333 | 0.968 |
| Baseline + Huffman | 244 | 6015.565 | 0.224 | 0.968 |
| Clustering | 892 | 5559.868 | 0.227 | 0.836 |
| Clustering + Pruning | 688 | 5721.021 | 0.230 | 1.017 |
| Clustering + Huffman | 240 | 4991.857 | 0.223 | 0.836 |
| Clustering + Pruning + Huffman | 128 | 5251.288 | 0.228 | 1.017 |
| †Clustering | 92 | 0.426 | 0.046 | 0.836 |
| **†Clustering + Quantization** | **60** | **0.271** | **0.043** | **0.834** |

- **18 x** reduction in model size

- **24,000 x** load time improvement

- **8 x** inference latency improvement

# Production Results

- One full night of ZTF alerts ~ 200K

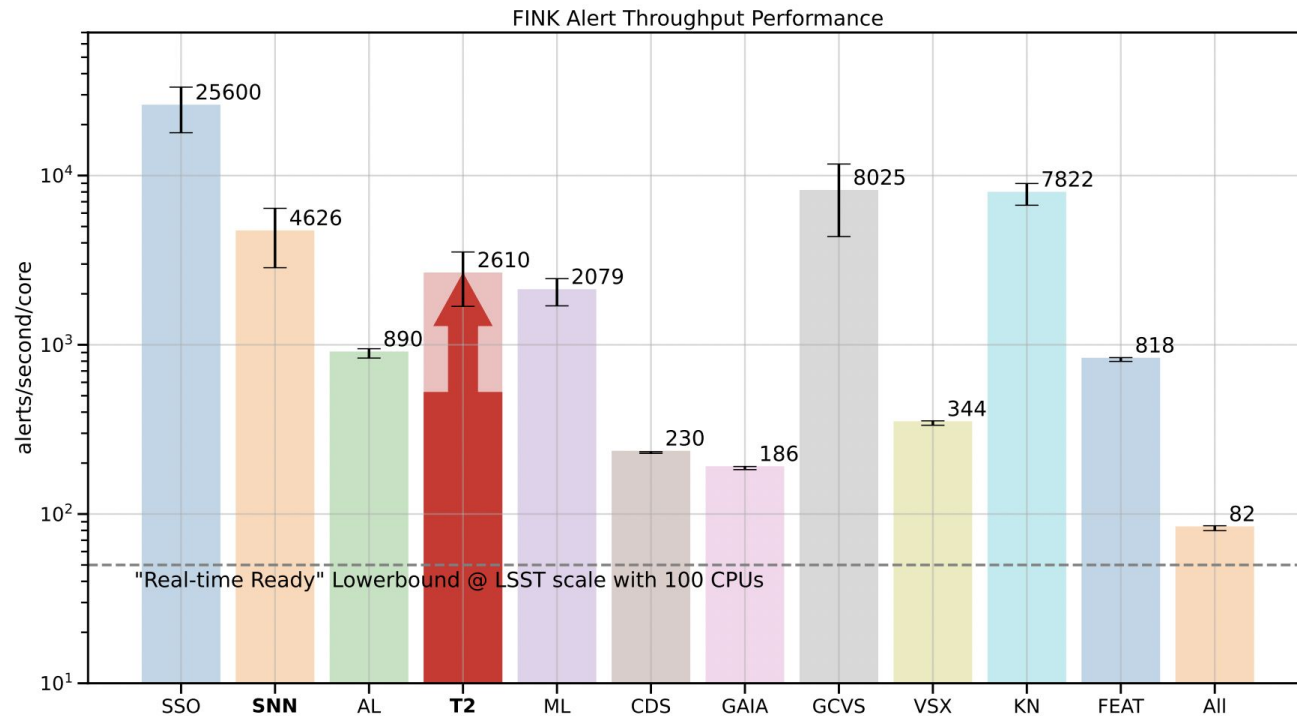- Require at least 2 points on light curve

- Averaged over 20 processing runs

# Production Results

- One full night of ZTF alerts ~ 200K

- Require at least 2 points on light curve

- Averaged over 20 processing runs

- **5x** throughput

FINK Alert Throughput Performance



alerts/second/core

- SSO: 25600
- SNN: 4626
- AL: 890
- T2: 2610
- ML: 2079
- CDS: 230
- GAIA: 186
- GCVS: 8025
- VSX: 344
- KN: 7822
- FEAT: 818
- All: 82

"Real-time Ready" Lowerbound @ LSST scale with 100 CPUs

# Takeaways...

- Use of deep compression can be your friend, *but be careful*

- Use FlatBuffers where possible



FINK Alert Throughput Performance

# Questions?