



### Machine Learning in the Multi-Messenger Era: Inference as a service and optimal light curve augmentation

Michael W. Coughlin University of Minnesota

On behalf of many MMA folks (A3D3, ZTF Machine Learning group)

November 30, 2021

# A Technical Ecosystem



Multi-Messenger science includes representation from many of the most interesting experiments today.





З





### So What's the Problem? Long Road from data to science



# Luckily, a lot of effort at this conference on problems of relevance!



# Luckily, a lot of effort at this conference on problems of relevance!





# Are we ready to apply ML in real time?

### Deep Learning Programs at Inference Time

#### Pros

- Robust modelling capabilities
- Real-time compatible

#### Cons

- Resource intensive
- May require frequent updates
- Effective use requires specialized knowledge, software, and often hardware



# Inference as a service



- Application for hosting trained networks and exposing them for inference via standardized client APIs
- Abstracts away details about models and their implementations
- Effectively leverages concurrent execution on heterogeneous computing resources
- Containerizations means portability and easy scale
- Centralized model repositories keep all users in sync

## Cleaning gravitational-wave data



11

### Detecting gravitational waves

- Gravitational-wave detection is basically a solved problem for Gaussian noise, which gravitational-wave data is not
- ML methods have been shown to have the capability of meeting the speed requirements for online searches, while also being more robust to data transients
- Only BBHs (short signals) so far







### LVK Inference as a Service: Deployment Scenarios



Online

- Latency sensitive, deploy locally to minimize data transfer time
- Using DeepClean to remove noise in real-time and make cleaned strain available to downstream analyses/ searches, including BBHNet

Offline

- Maximize throughput to minimize time to completion (subject to cost constraints)
- Cloud resources leverage economies of scale
- Cleaning one month of O3 data with DeepClean
- End-to-end ensemble with DeepClean and BBHnet to estimate event likelihood over ~27 hrs. of O2 data



### Inference as a Service: for Streaming Data







(c)

- DeepClean and BBHnet perform inference on fixed-length snapshots of time series
  - Rate at which snapshots are sampled fixed by an inferencetime parameter *r*, the inference sampling rate
- High values of r compared to the length of the frame lead to substantial data overlap and redundant data transfer from client to server
- Host a "snapshotter" model on the server that maintains the current snapshot as a state
  - Only stream state updates
  - Updated snapshot gets passed to downstream models
  - Introduces a potential sequential bottleneck
- DeepClean also has overlapping output data
  - Aggregating between overlap incurs extra latency
  - Currently adopting "fully-online" solution

Gunny et al. 2021: 2108.12430

(b)

#### Performance vs aggregation latency





UMN

aggregation latency the amount of data (seconds) to be excluded from the end of the segment due to quality degradation.



## Offline Use Cases



- Offline DeepClean shows the advantages of switching to laaS model: CPU-only node with ~10x reduction in processing time, adding GPUs gives another ~5x
- Ensemble model shows economies
  of scale of laaS paradigm
- Processing time decreases linearly with # of nodes, cost stays constant
  - Optimal point at "infinite" scale
  - Scale achieved with minimal additional engineering overhead
- Ensemble leverages multiple framework backends, all invisible to client users



#### Gunny et al. 2021: 2108.12430

# Conline Use Cases



- Concurrent execution of laaS model important at lower inference rates, keeps demand for scale low
- Scale becomes more important at higher frequency inference rates
- Bottleneck is currently sequential update to the input "state"
- Optimizing this step via HPC unlocks more advantage from additional GPUs

17



- Use HEPCloud framework to run larger tests with multiple clients/servers
  - Tests use cloud
    resources (GCP),
    submitted through
    HTCondor
- Jobs synchronized to start all at once, mimic realistic environment
- Able to sustain processing for full length of job
- Provides a means to manage/run large amount of jobs





# How can we support and perform better observations?

## The Observational Landscape



## Landscape of Optical TDA

![](_page_20_Figure_1.jpeg)

- The night sky is imaged at 17.5 mag by ASAS-SN (both hemispheres)
- The northern sky is covered by ATLAS, ZTF, and PS-1 to 19, 20.5, 21.5 over roughly two nights (ZTF issues real time, data-rich alerts)
- BlackGEM (21-22 mag; Chile) will start routine operation within this year
- Rubin is expected to become operational in 3 to 4 years

## Two Different Approaches

![](_page_21_Figure_1.jpeg)

- Photometric detection followed by spectroscopic classification
- Possible for surveys which are shallow (20 mag or so)
- Spectral classification can be undertaken by existing telescopes
- Photometric detection followed by multi-band time series
- Large samples of faint objects –Much of the analysis will be statistical –Use clever techniques and filter out a small subset for further follow up

![](_page_22_Picture_0.jpeg)

![](_page_22_Picture_1.jpeg)

P48 Discovery

P200 Spectroscopy

![](_page_22_Picture_4.jpeg)

## The Technical Landscape

![](_page_23_Figure_1.jpeg)

## The Technical Landscape

![](_page_24_Figure_1.jpeg)

## Wide Field Follow-up

![](_page_25_Figure_1.jpeg)

### **Observing Scenarios**

![](_page_26_Figure_1.jpeg)

### What data do we need?

- Often, a photometric light curve is all you have available to classify it.
- Due to the many follow-up systems we have available, desire to design a system that optimizes the differentiation between models for kilonovae and other fast transients.
- Can use ML methods to speed up inference on each potential counterpart object, including when performing the GW and EM inference.

![](_page_27_Figure_4.jpeg)

![](_page_27_Figure_5.jpeg)

### Transient Filtering

![](_page_28_Figure_1.jpeg)

#### [Andreoni, Coughlin+2021, 2104.06352]

### Why do we want fancy strategies? Kilonovae - Hard to find

![](_page_29_Figure_1.jpeg)

## Landscape of Optical TDA

![](_page_30_Figure_1.jpeg)

- The night sky is imaged at 17.5 mag by ASAS-SN (both hemispheres)
- The northern sky is covered by ATLAS, ZTF, and PS-1 to 19, 20.5, 21.5 over roughly two nights (ZTF issues real time, data-rich alerts)
- BlackGEM (21-22 mag; Chile) will start routine operation within this year
- I SST is avanceed to become energy on a stad years

31

## Two Different Approaches

![](_page_31_Figure_1.jpeg)

- Photometric detection followed by spectroscopic classification
- Possible for surveys which are shallow (20 mag or so)
- Spectral classification can be undertaken by existing telescopes
- Photometric detection followed by multi-band time series
- Large samples of faint objects –Much of the analysis will be statistical –Use clever techniques and filter out a small subset for further follow up

## Value-driven Real-time follow-up

- As is, data will be insufficient for full science inference without additional follow-up
  - •e.g. extracting physics from light curves
- Need to perform value-driven follow-up
  - Volume of alerts far exceeds the ability to follow-up with limited and/or expensive follow-up resources
- Augment sparse Rubin LCs to improve constraints on SALT2 (supernova) models
- We augment photometry to (branch-normal) SN Ia LCs from ZTF-I public survey (g and r) using P48 in g,r, and i
  - i-band important for precisely estimating H0 (Burns+ 2018)
  - Second peak could help probe SN Ia explosion mechanisms
- Broadly an optimal real-time resource allocation problem and not restricted to SALT2

![](_page_32_Figure_10.jpeg)

![](_page_32_Picture_11.jpeg)

![](_page_32_Picture_12.jpeg)

Niharika (Ari) Sravan

33

![](_page_33_Picture_0.jpeg)

10-20 % Median improvement in parameters over random allocation

Gap filling Resolves phase with high variability/ diversity:

Around peaks and valleys

![](_page_33_Figure_4.jpeg)

![](_page_34_Picture_0.jpeg)

Interesting notes:

More in g due to sparser sampling

3-5% more improvement for SNe Ia > 18.5mag

Even better prospects for Rubin

![](_page_34_Figure_5.jpeg)

![](_page_34_Figure_6.jpeg)

![](_page_34_Figure_7.jpeg)

# So... what then?

#### NMMA: A Fully Bayesian Joint-Inference Pipeline

- gravitational-wave data analysis using parallel bilby
- kilonova modelling with various models (Bulla, Kasen, etc.)
- gamma-ray burst afterglow fits (also supernova models from sncosmo)
- chiral effective field theory to simulate the neutron-star EOS
- neutron-star maximum mass and NICER constraints, fits to relate ejecta parameters to progenitor parameters using numerical relativity

Used for many analyses at this point: Ahumada et al. 2105.05067, Dietrich et al. 2002.11355, Tews et al. 2007.06057, Pang et al. 2105.08688, Huth et al. 2107.06229

![](_page_35_Figure_8.jpeg)

# NMMA: A Fully Bayesian Joint-Inference Pipeline

![](_page_36_Figure_1.jpeg)

Peter Pang

NIKHEF

- Extract science (and filtering criteria) in one place!
- For example, immediately extract information on ejecta, neutron star physics, and cosmology (in case of a host galaxy).

# Online Filtering

Ø

![](_page_37_Figure_1.jpeg)

 $\equiv$  All unreads

Threads

റി All DMs

: More

# decam

gattini

#

#

# git

#

#

#

# nmma

# nmma-bot

# operations-research

grb

#### # nmma-bot ~

Today ~

![](_page_37_Picture_4.jpeg)

UMN

![](_page_37_Picture_6.jpeg)

![](_page_37_Figure_7.jpeg)

growth-mma-restbot APP 11:42 AM Hi ztfrest! You are interested in ztfrest fitting, right? Let me get right on that for you. Name: ZTF21abneypf @ Mentions & reactions Model: Bu2019lm log(Bayes): -30.204316408615664 ℅ Slack Connect log(Evidence): -30.204316408615664 ± 0.10116267044829176 Bu2019Imcorner ZTF21abneypf - Channels A ampel\_mm # gemini general grb200826a # grb201130a # grb210510a 11:43 Bu2019Imlightcurve ZTF21abneypf • # grb210529b # ic201021a # ic210210a # ic210629a neutrinos

![](_page_37_Figure_9.jpeg)

![](_page_37_Figure_10.jpeg)

#### Message #nmma-bot

#### [Barna, Reed, et al., in prep]

#### 38

### Improving community follow-up

![](_page_38_Figure_1.jpeg)

![](_page_39_Picture_0.jpeg)

![](_page_39_Picture_1.jpeg)

#### Join us!

- Monday 9 am Central: Technical MMA call (anything related to gravitational-wave counterpart searches)
- Tuesday 2 pm Central: A3D3 ML Detection meeting
- Thursday 8:30 am Central: A3D3 KAGRA meeting
- Friday 2 pm Central: A3D3 Inference as a Service meeting

![](_page_40_Picture_0.jpeg)

### **Summary and Perspectives**

#### IAAS

- IaaS model represents a powerful way to bring advantages of deep learning to bear in gravitationalwave astronomy
- Deploying and optimizing laaS pipelines requires aligning benefits of scale with constraints of problem
- Several applications currently under development to increase number and speed of event detections during O4

#### ML based follow-up

- Technique interpretable
  - Good start before invoking reinforcement learning (also for benchmarking)
- Data acquisition failure tolerant
  - Latency intolerant
- Main uncertainty is the reliability of simulated augmented photometry
  - Need verification with real data
- Extensions of work include:
  - Variable observing cost, observing season based budget
  - Other choices of utility including prior building, model discrimination

My own perspective on the areas of greatest need:

- Can the initial promise of ML applications in terms of detection and PE make its way from BBH signals (short) to BNS signals (long)?
- Can ML provide optimal follow-up strategies to rule in and out specific transients as sources given limited telescope time and sensitivity?
- Is ML the key to a truly MMA pipeline, with inference on GW, optical, GRB, etc. data sets?

![](_page_41_Picture_0.jpeg)

### Thank you!

### MMA Equation of State Constraints

![](_page_42_Figure_1.jpeg)

![](_page_43_Picture_0.jpeg)

#### The Goal

- Given an observing budget decide how to augment LCs in real-time (adapt to collected data) such that we maximize expected utility (EU) for the augmented LC (realized at the end of the episode)
  - Utility conveys our preference for an outcome given a decision. We choose utility as pseudo\* A-optimality = minimum SALT2 parameter variances in sncosmo
  - Since we do not assume that redshift is known as the SN is taking place, we solve for it and minimize its uncertainty as well.
     Assume known SN sky location, MW extinction.

#### The Algorithm

- On each day estimate EU of action space {no action, g, r, i, gr, ri, ig, gri} given observed data and expected data (stochastic) under no action.
   Remaining budget allocated randomly\*\*
  - Outcome states given actions and expected future data estimated using encoder-decoder
     LSTM trained on 105 simulated ZTF SNe Ia
- Take modal action with least cost having max EU\*\*\* per N simulated future outcomes
  - Augmentations from 2-D Gaussian Process fit to full LC and fed back for the next day

\* because sncosmo solves chi^2 minimization max likelihood
 \*\* substitutes expected **optimal** actions for expected **naïve** actions
 \*\*\* -greedy and tuned to maximize median A-optimality for all validation SNe Ia

#### Sraven et al. 2021: submitted