# Finding exoplanets with TESS & AI
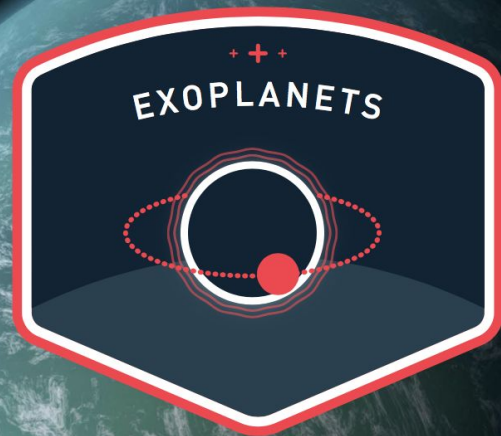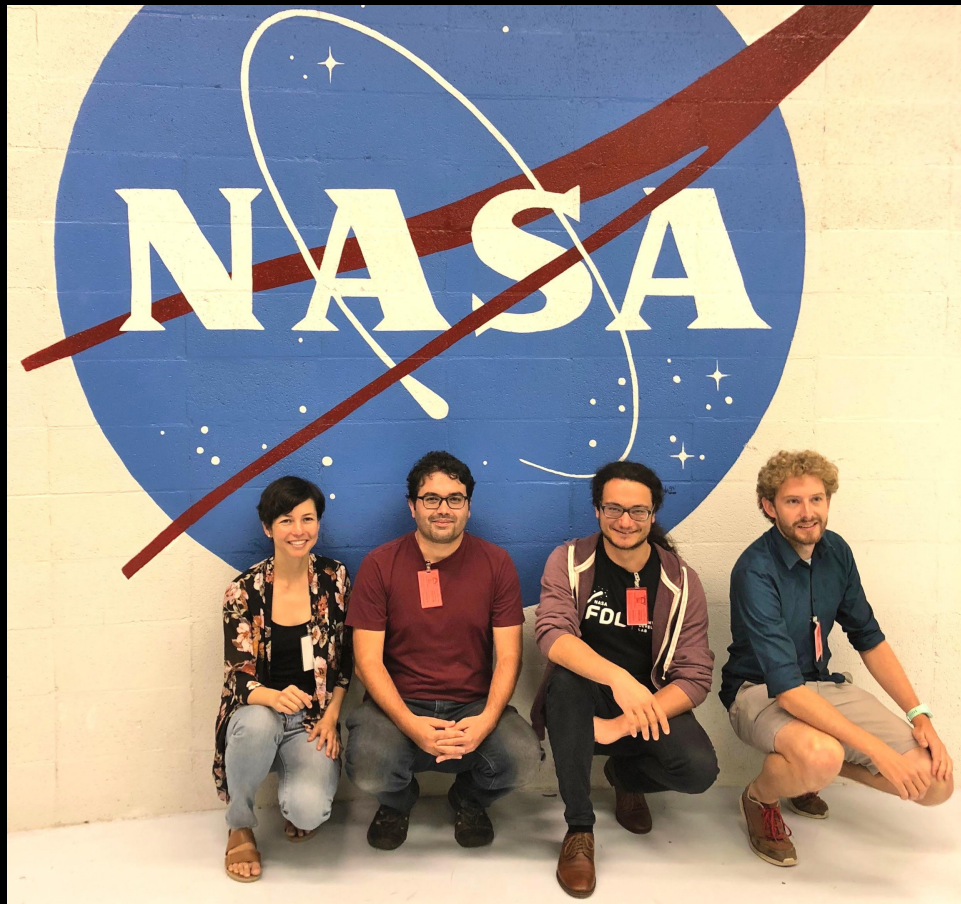
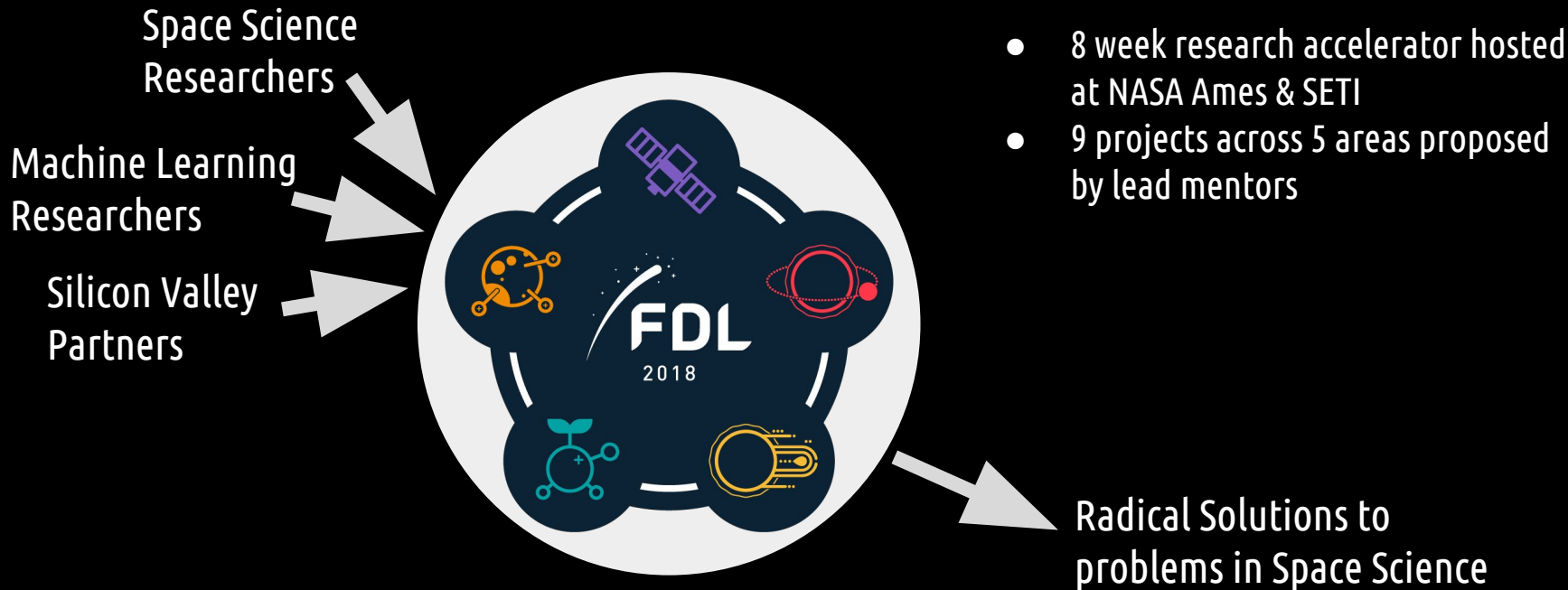**Megan Ansdell**, Yani Ioannou, **Hugh Osborn**, Michele Sasdelli,

+ Jeff Smith, Jon Jenkins, Doug Caldwell, Adam Lesnikowski, Chedy Raissi, Massimo Mascaro
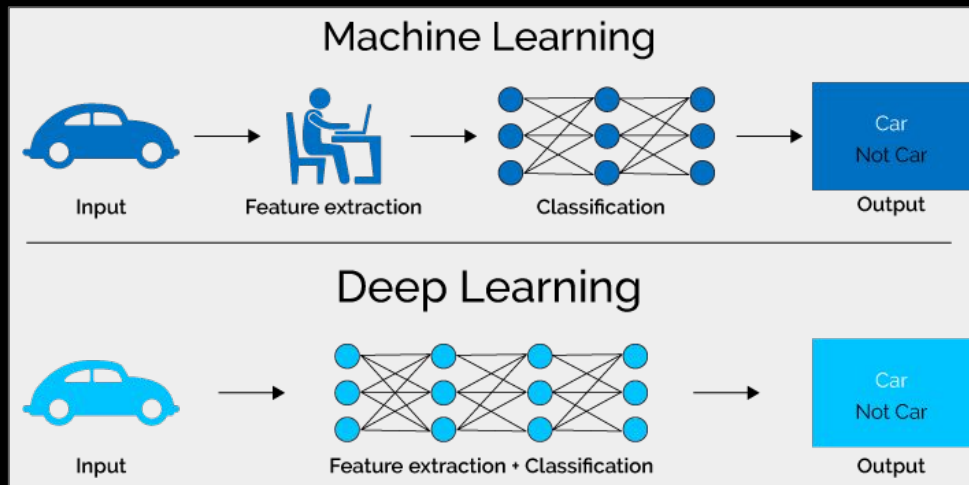
# The Team



Megan Ansdell, Yani Ioannou, Michele Sasdelli, Hugh Osborn

FDL

# Frontier Development Lab

Space Science
Researchers

Machine Learning
Researchers

Silicon Valley
Partners

FDL
2018

- 8 week research accelerator hosted at NASA Ames & SETI
- 9 projects across 5 areas proposed by lead mentors

Radical Solutions to problems in Space Science

FDL

SETI INSTITUTE · intel AI · SR SPACE RESOURCES.LU · XPRIZE · Google Cloud · nVIDIA · LOCKHEED MARTIN · kx · IBM · KBRwyle We Deliver
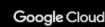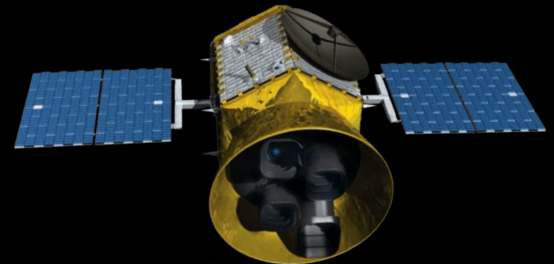
# Deep Learning

- **Excellent at classification problems when:**
  - Lots (N> 10 000s) of labelled data
  - Signal is complex to model (not True for "planet" but true for non-planet!)

- **Once trained, a deep learning algorithm is:**
  - Far faster than classical at performing classification

- **But:**
  - Large computing infrastructure to train
  - Must set aside much of the data as a training set



Movie credit: Denis Dmitriev

# Deep Convolutional Neural Networks

- **Machine Learning (ML)**:
  *models learn features from data*

- **Deep Learning**:
  *layers build increasingly complex features*

- **Neural Network (NN)**:
  *model learns weights of nodes*

- **Convolutional Neural Network (CNN)**:
  *exploits spatial structure in data*

- **Binary Classification**:
  *final layer outputs single number from 0-1*

Movie credit: Denis Dmitriev

# The Data - TESS

# TESS "Data" (TSOP-301 simulations)

Target Pixel Files (simulated planets)

16k stars per sector

2min cadence

TPFs

data:

3 simulated sectors

Real data (from 2019):

One sector per month

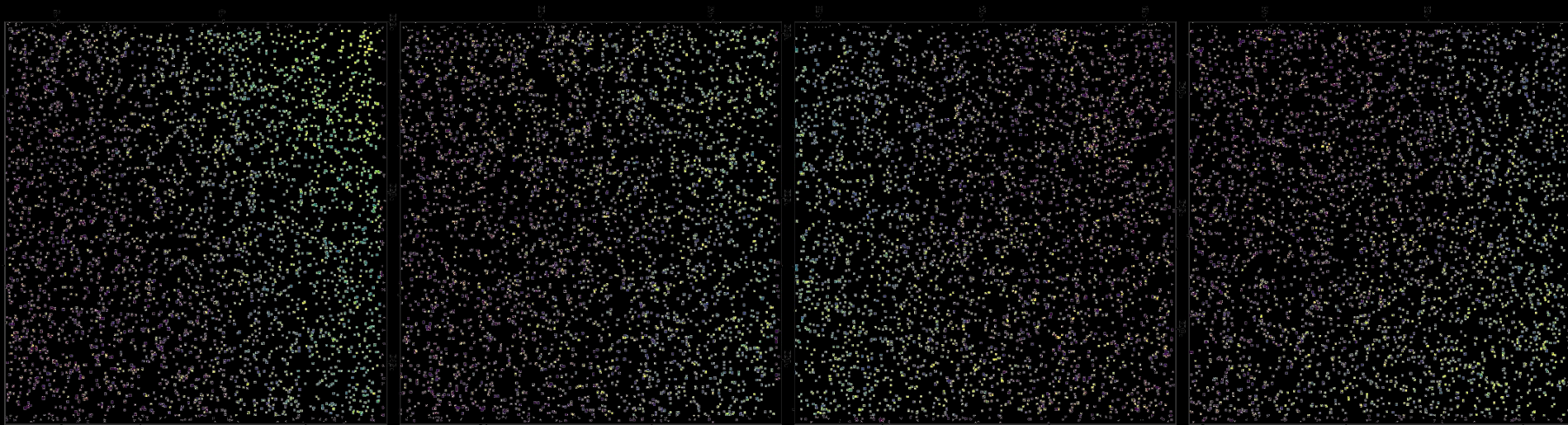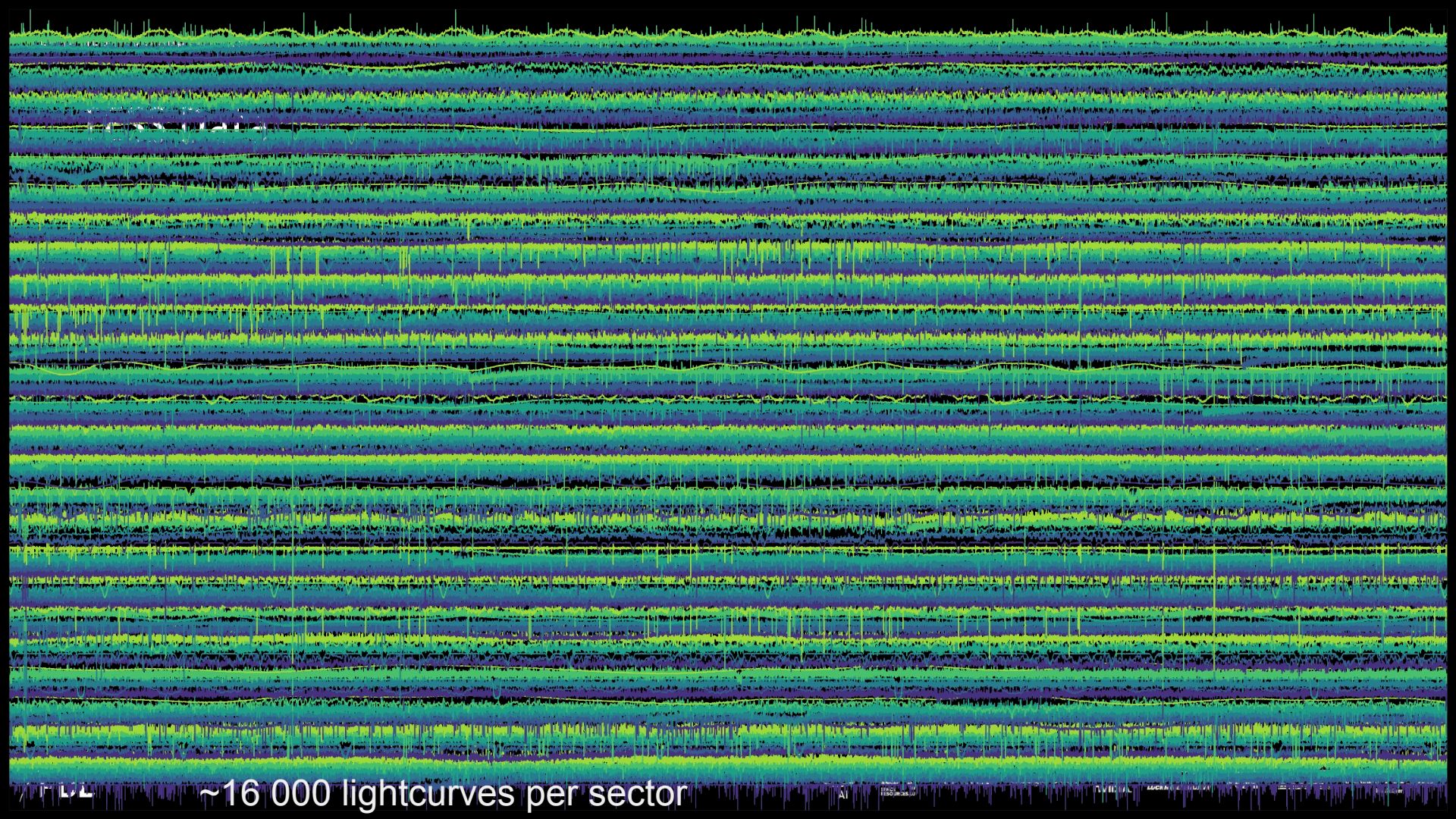Require rapid, accurate, planet identification

# The Data - TESS



Per month:          4 CCDs          ~16 000 target stars          21 000 images

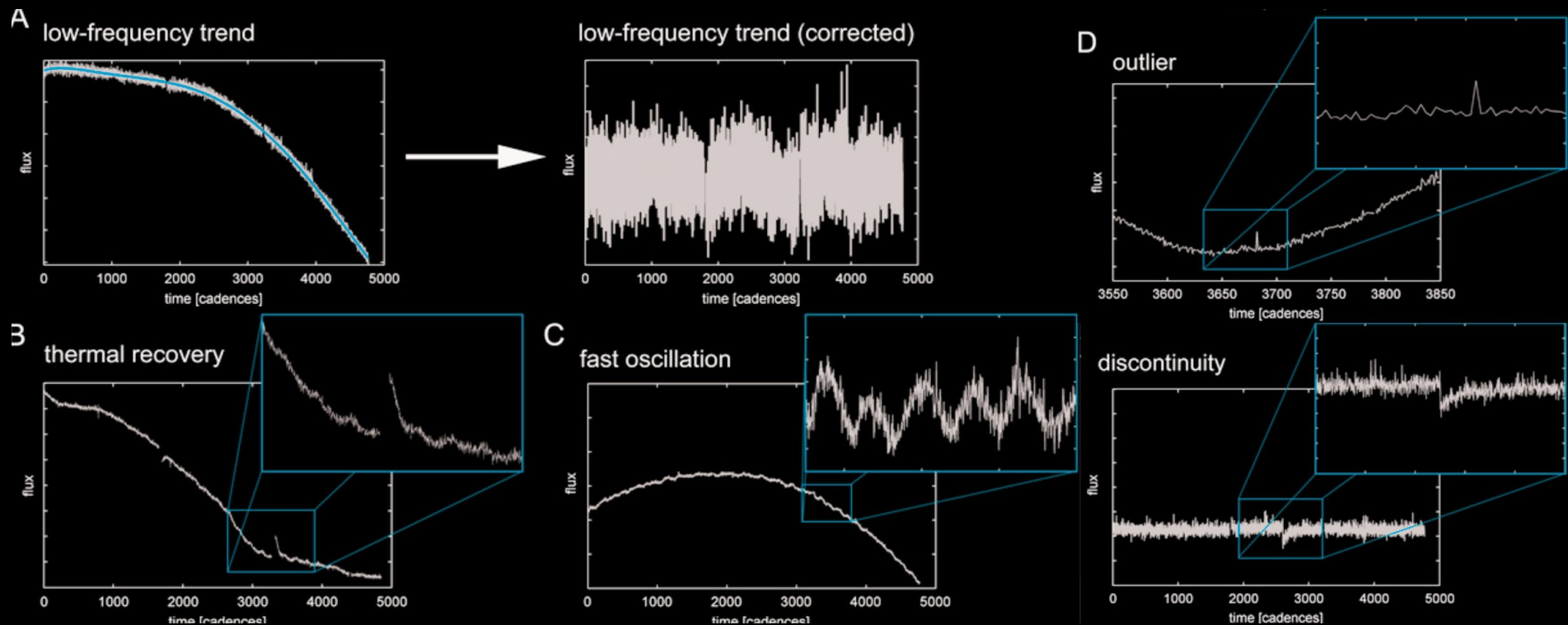Our dataset: TSOP-301. 4 simulated sectors

# TESS - the data bottleneck



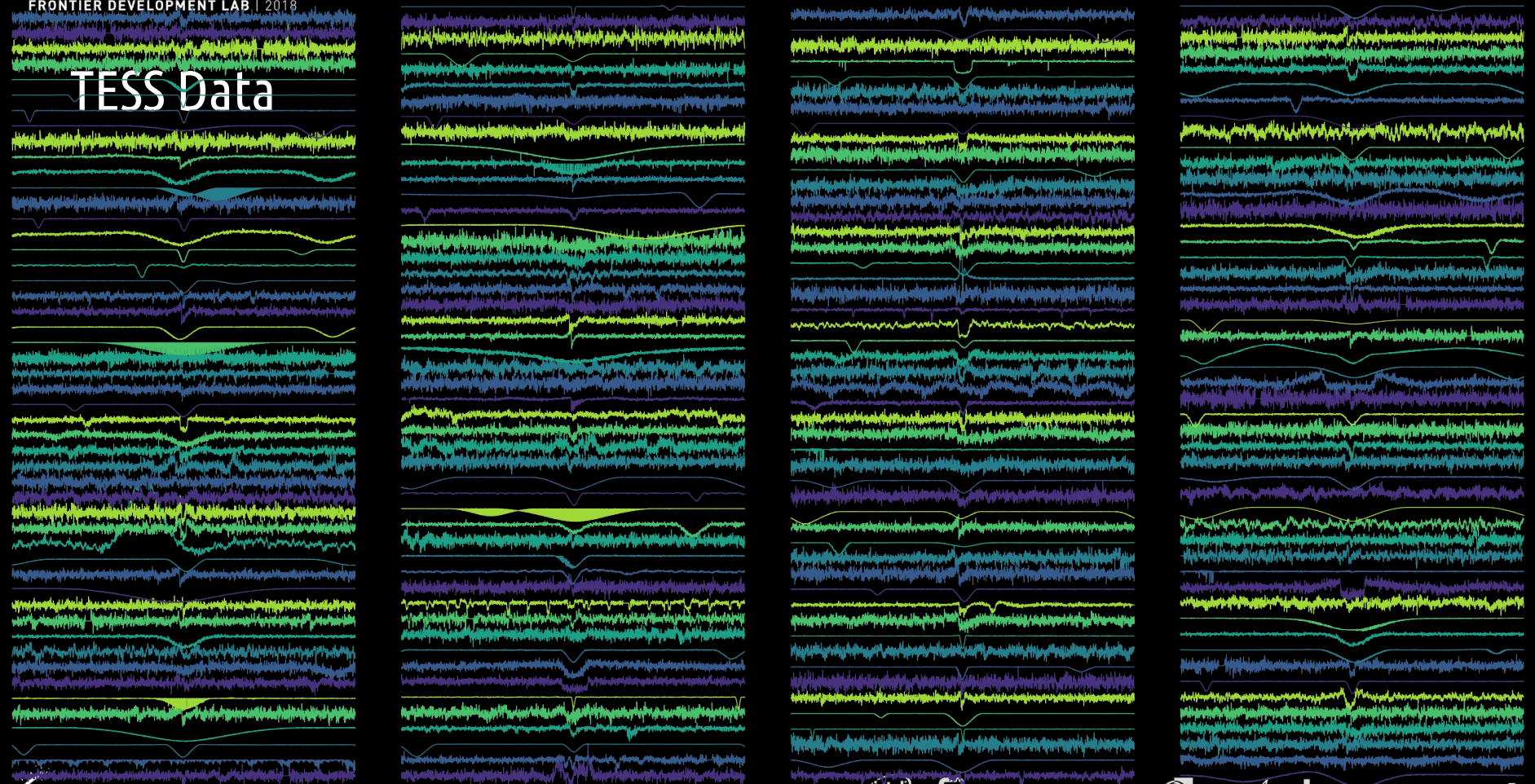4 simulated sectors     ~16 000 stars per month     ~4 000 hours of video

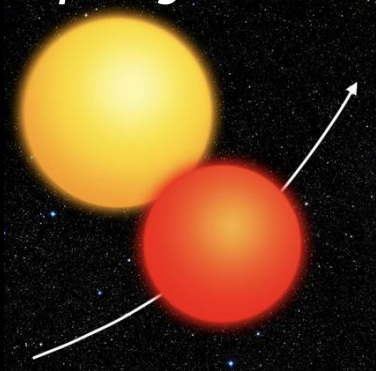~16 000 lightcurves per sector

# TESS Systematics

TESS Data

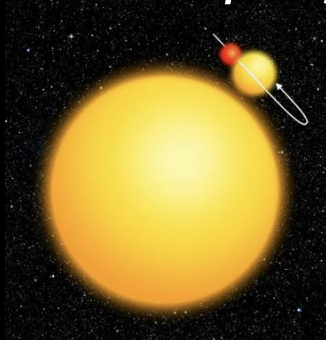~5000 candidates per sector

FDL

TESS Data

FDL

TESS Data

Planet

Not Planet

FDL

intel AI    SK SPACE RESOURCES.LU    XPRIZE    Google Cloud    nVIDIA    LOCKHEED MARTIN    kx    IBM    KBRwyle We Deliver

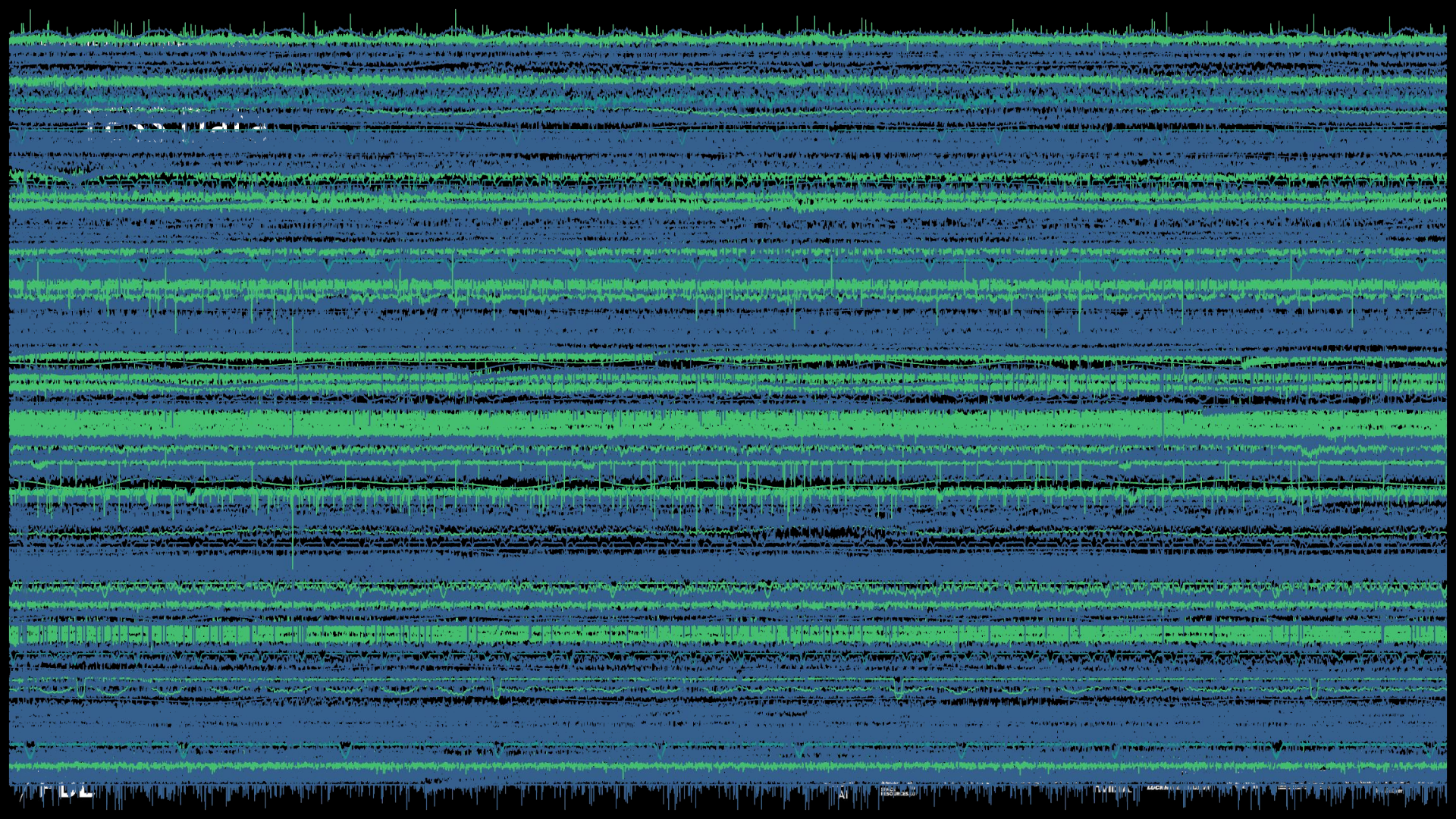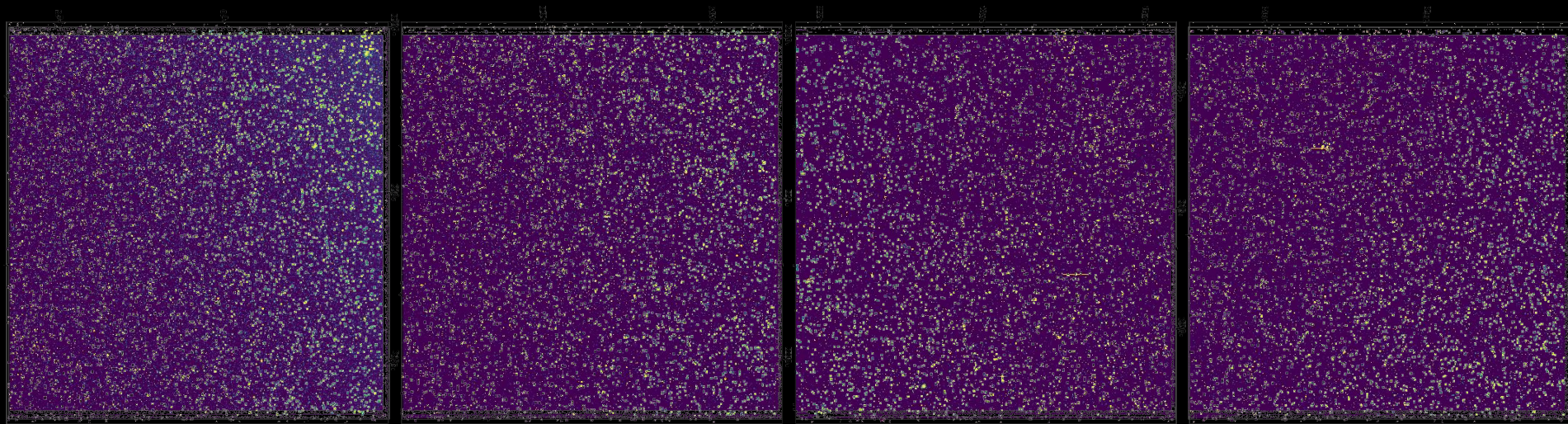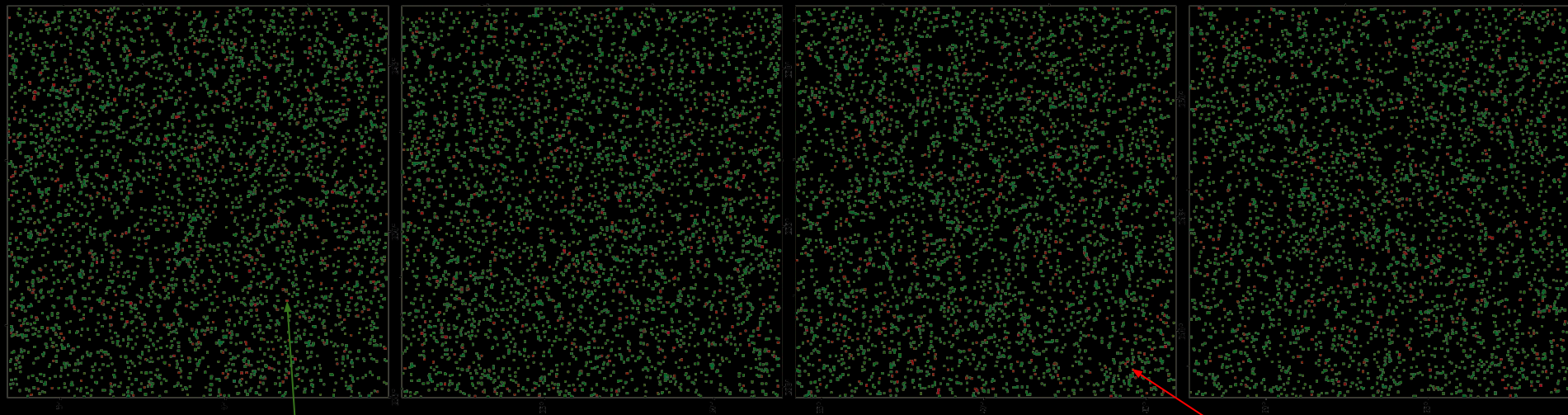# TESS Data

# TESS Data



Not Planet

Planet

# Current Classification Technique

- Statistical/automated methods are used to whittle down candidates

- Manual vetting is still common

- Team of 18 humans: 94% accuracy
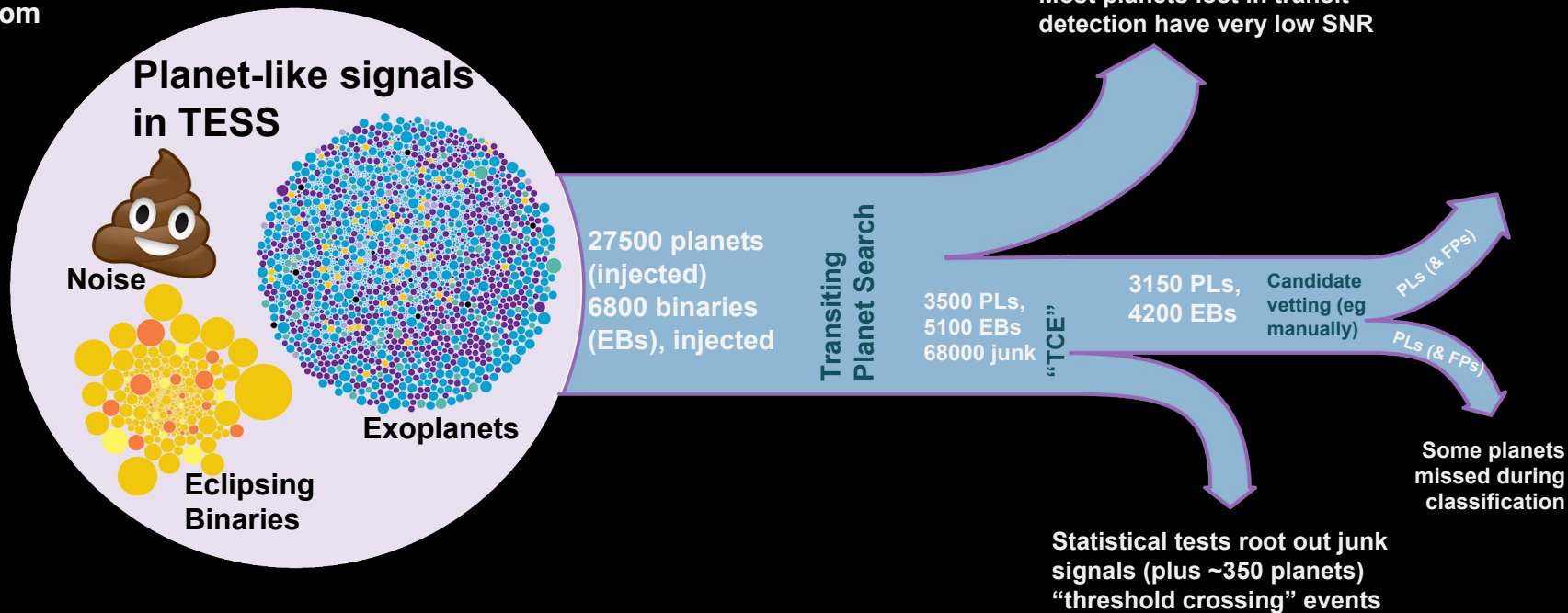
- ~300 human hours per sector

The telescopes are waiting...

# From the detections...
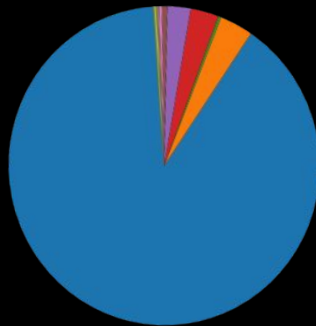
**Classical planet search in TESS**
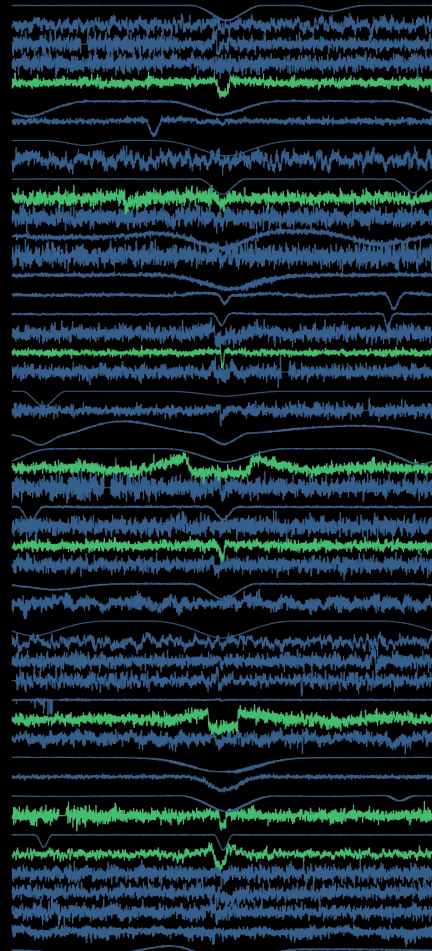
**Numbers from TSOP-301 simulation**

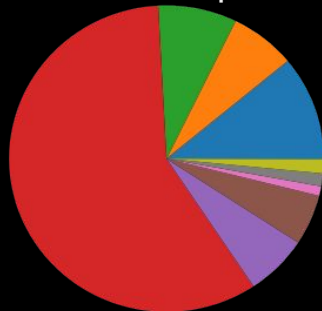**Planet-like signals in TESS**

Noise

Eclipsing Binaries

Exoplanets

27500 planets (injected)
6800 binaries (EBs), injected

Transiting Planet Search

3500 PLs, 5100 EBs 68000 junk

"TCE"

3150 PLs, 4200 EBs

Candidate vetting (eg manually)

PLs (& FPs)

PLs (& FPs)

Most planets lost in transit detection have very low SNR

Some planets missed during classification

Statistical tests root out junk signals (plus ~350 planets) "threshold crossing" events

FDL

intel AI    SR SPACE RESOURCES.LU    XPRIZE    Google Cloud    nVIDIA    LOCKHEED MARTIN    kx    IBM    KBRwyle We Deliver
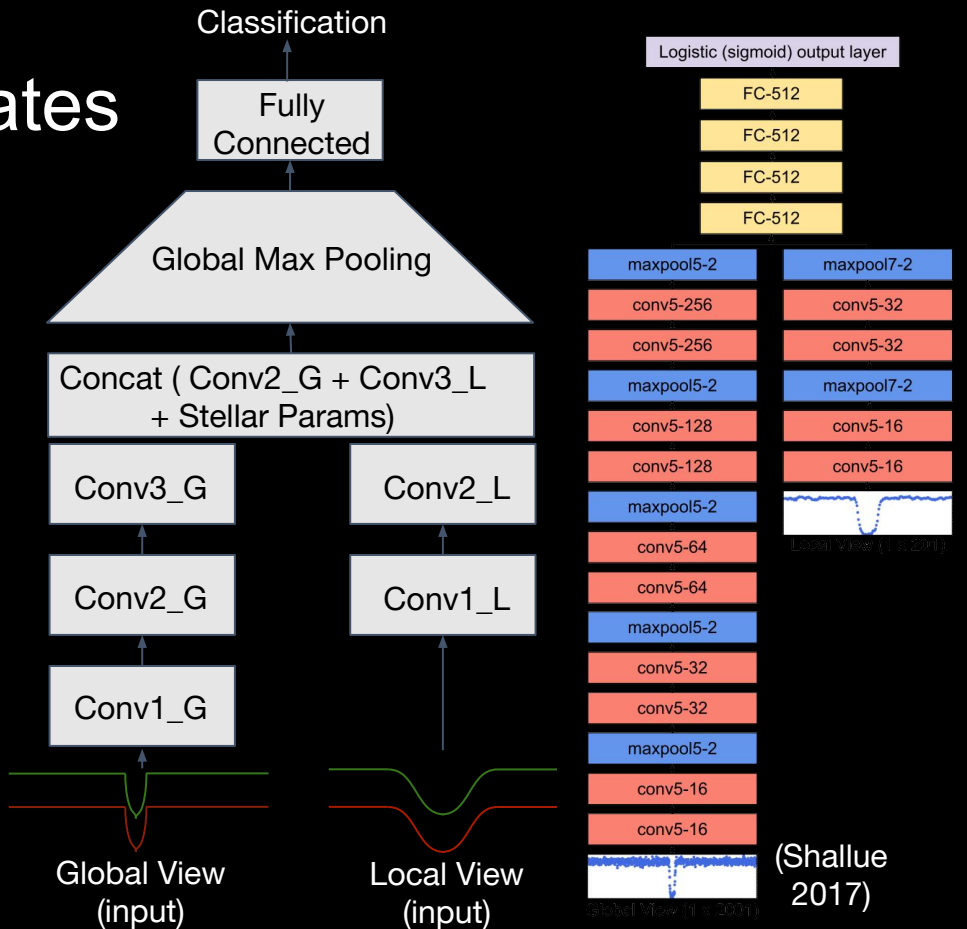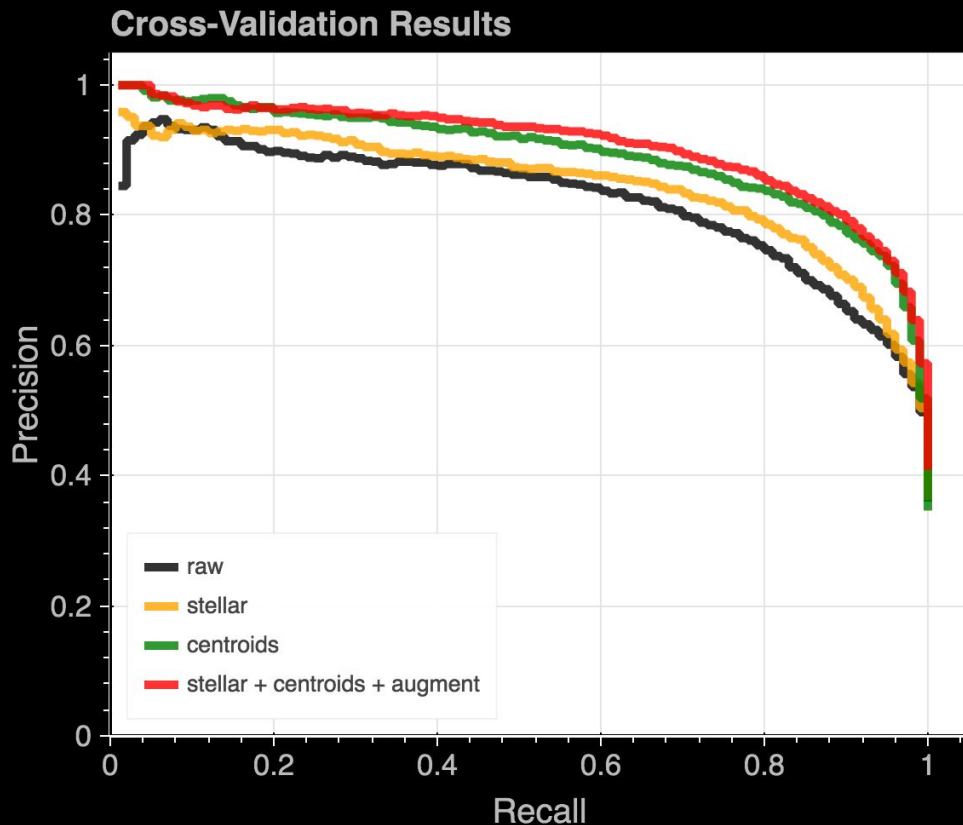
# Classifying TESS Candidates

Adapted Shallue et al

- Made model smaller (0.06%)
- Added stellar parameters
- Added motion of star (centroid)
- Mini-batch balancing to account for label imbalance

Classification

Fully Connected

Global Max Pooling

Concat ( Conv2_G + Conv3_L + Stellar Params)

Conv3_G | Conv2_L

Conv2_G | Conv1_L

Conv1_G

Global View (input) | Local View (input)

Logistic (sigmoid) output layer

FC-512
FC-512
FC-512
FC-512

maxpool5-2 | maxpool7-2
conv5-256 | conv5-32
conv5-256 | conv5-32
maxpool5-2 | maxpool7-2
conv5-128 | conv5-16
conv5-128 | conv5-16
maxpool5-2 | Local View (input 201)
conv5-64 |
conv5-64 |
maxpool5-2 |
conv5-32 |
conv5-32 |
maxpool5-2 |
conv5-16 |
conv5-16 | (Shallue 2017)
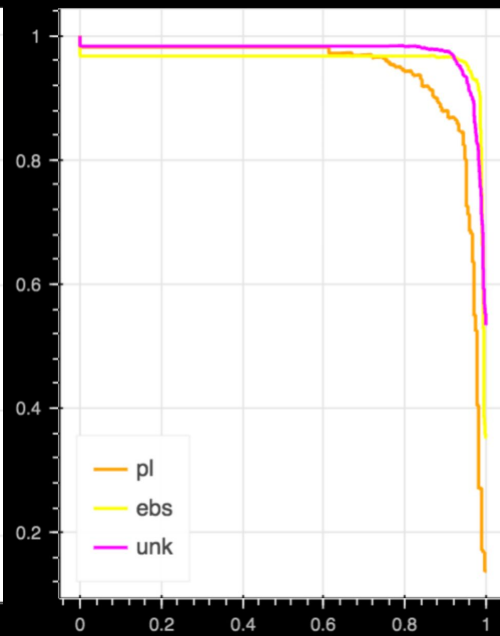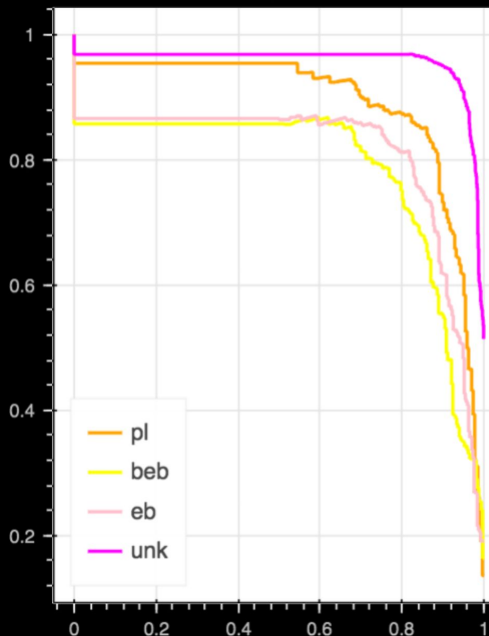
Global View (input 2001)

# From the detections...

- Average precision on Kepler ~96%

- Recovers ~90 more planets than

  Shallue et al (on single model)

- Model is ~500 times smaller

- Similar precision on TESS data

- Can be run in minutes not hours!

**Cross-Validation Results**
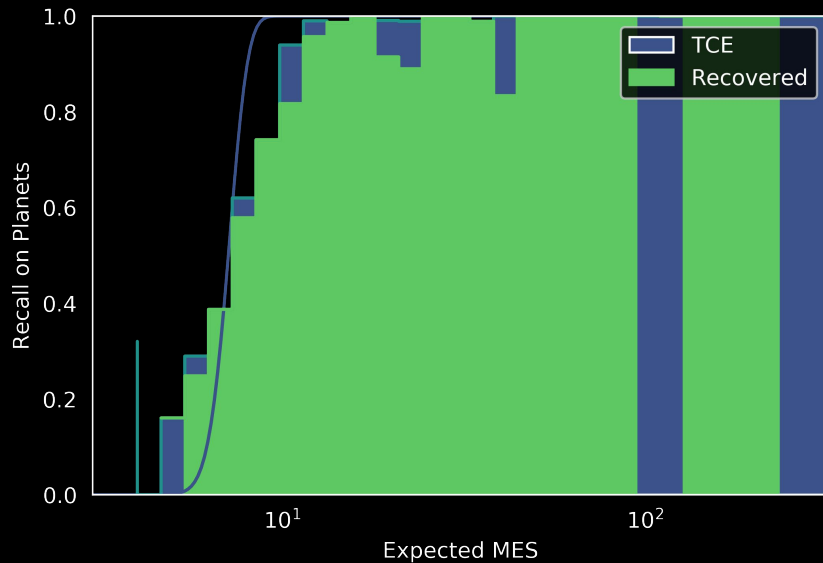


FDL

# From the detections...

- Also developed multi-class models on TESS data.
- Useful for follow-up!
- 3-class > 4-class
- Slightly lower precision on planets.
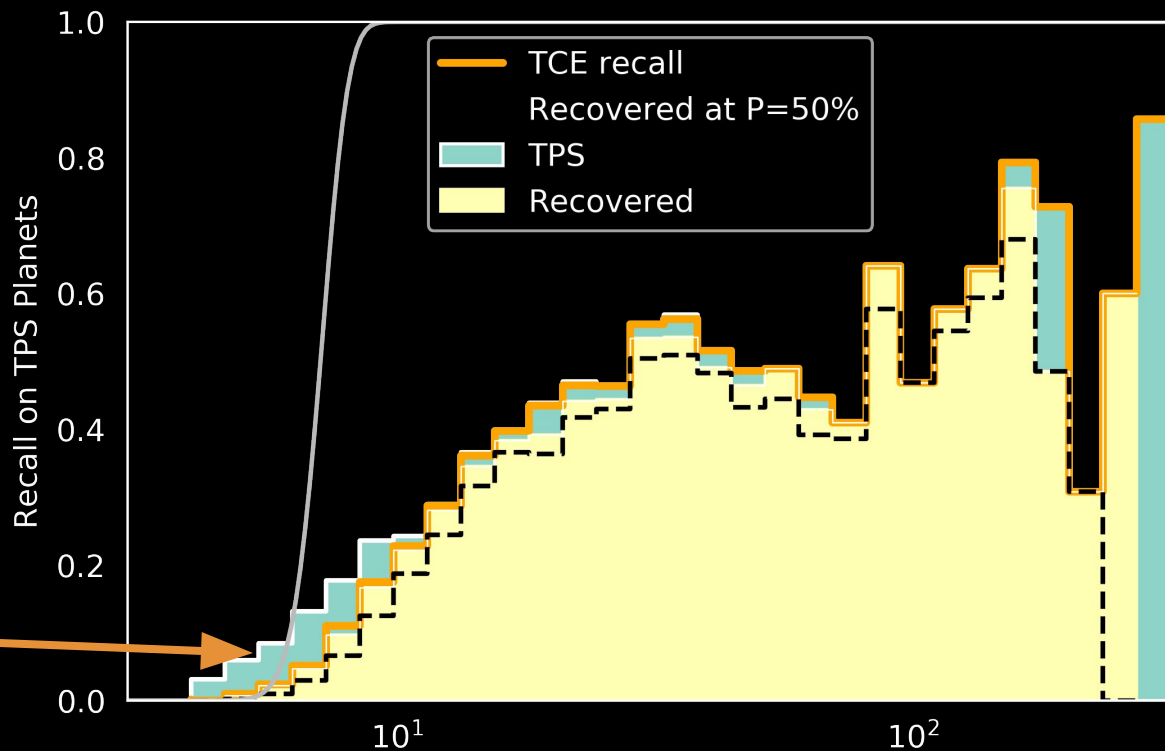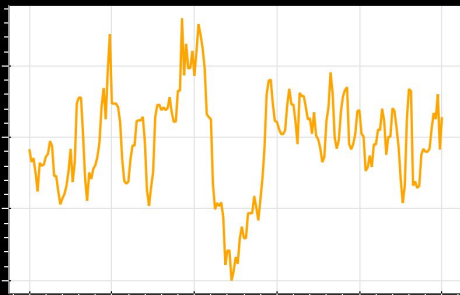
# From the detections...

Kepler vs Shallue

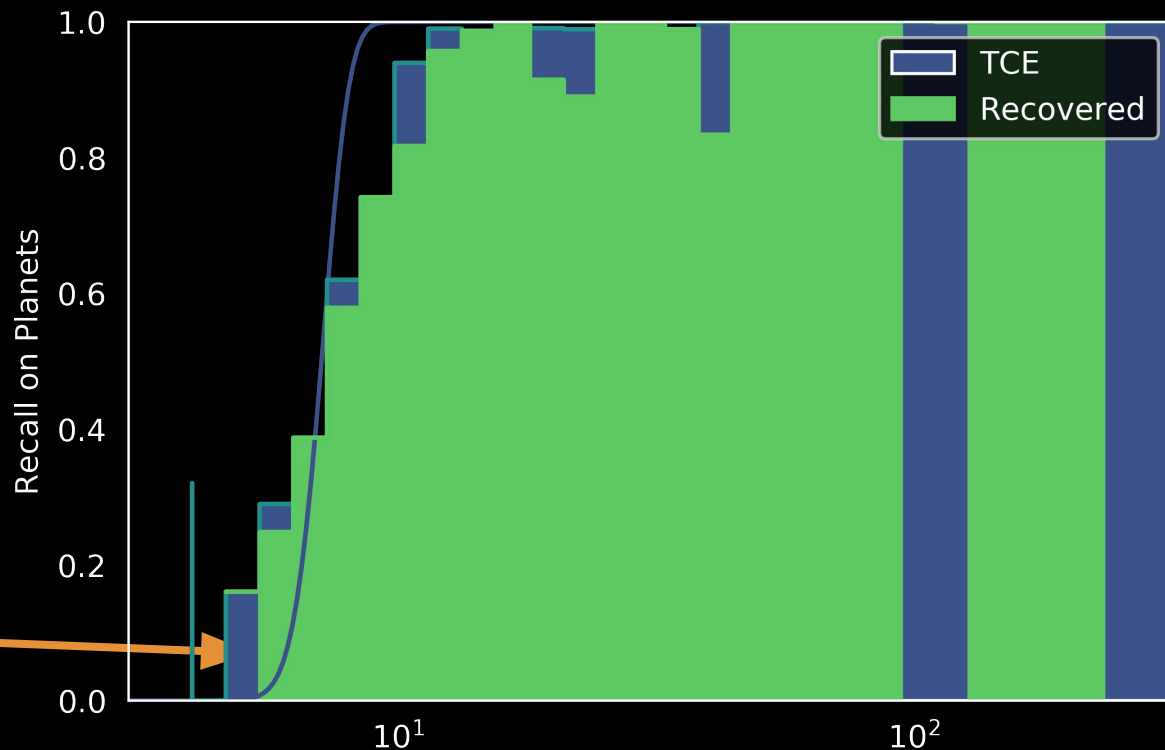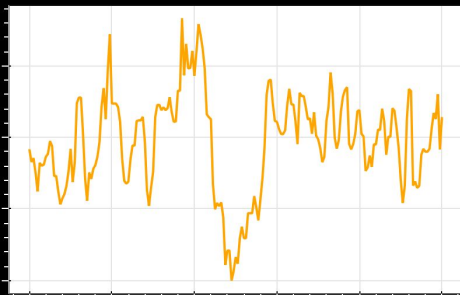Up to 650 more planets recovered (at precision of 0.9 on single model)

# From the detections...

- Can replicate the classical statistical threshold tests
- Recovered planetary signals missed by the classical tests

# From the detections...

- Can replicate the classical statistical threshold tests
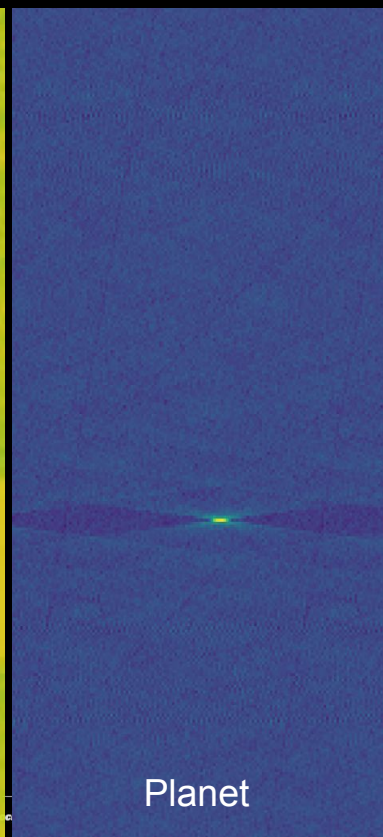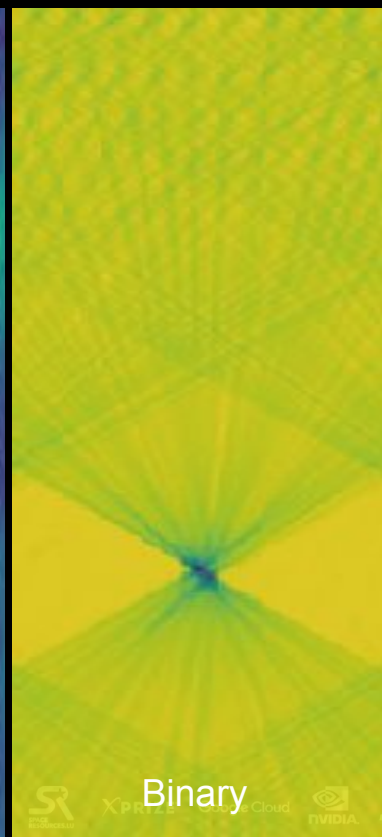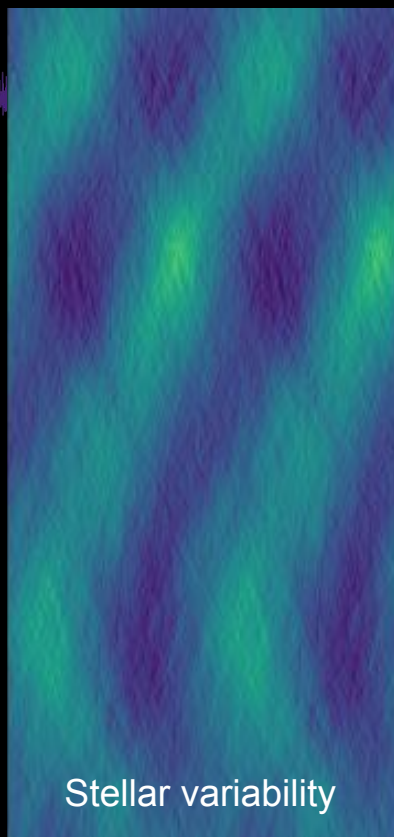- Recovered planetary signals missed by the classical tests
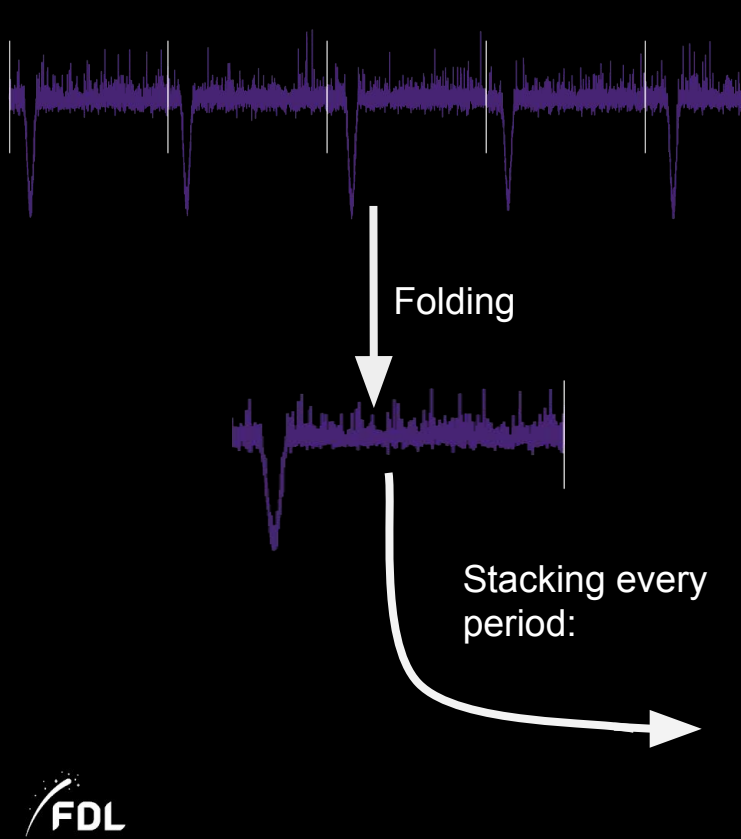
# Future - applying to real TESS data

- **Directly transfer network to real flight data…**
  - Noisy data
  - Is the data (& systematics) the same as in training?

- **Use human-vetted labels for TESS candidates…**
  - Noisy labels
  - Slow
  - Exactly what we're trying to replace!

- **Inject known signals into real flight data**
  - Can train network rapidly with real data and noise.
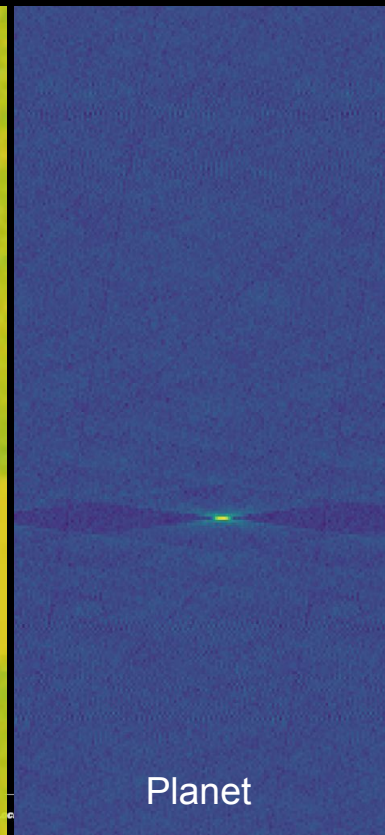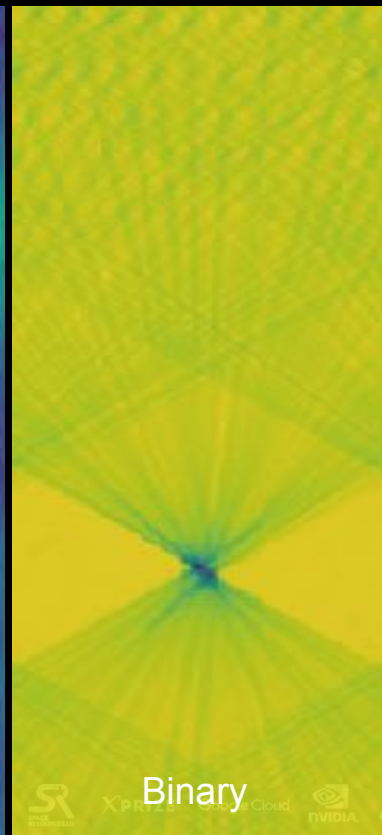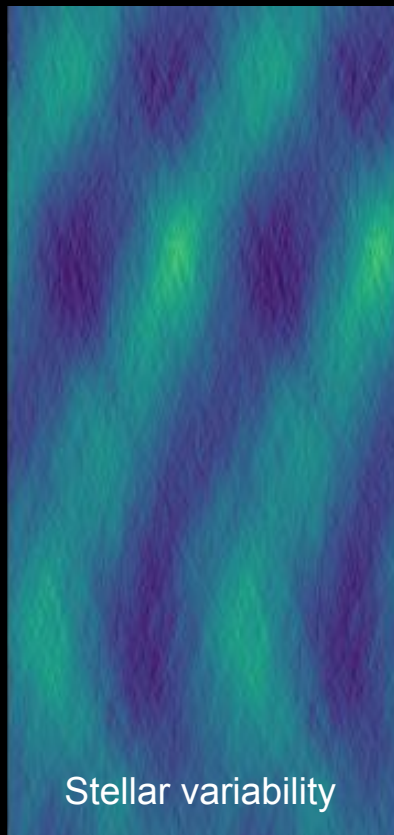  - Simulations and network re-training can be part of pipeline.

# From the lightcurve...



Folding

Stacking every period:

Stellar variability

Binary

Planet

# From the lightcurve...

- Applies ResNet50
- 91% accuracy on Kepler candidates
- Promising but need more time.



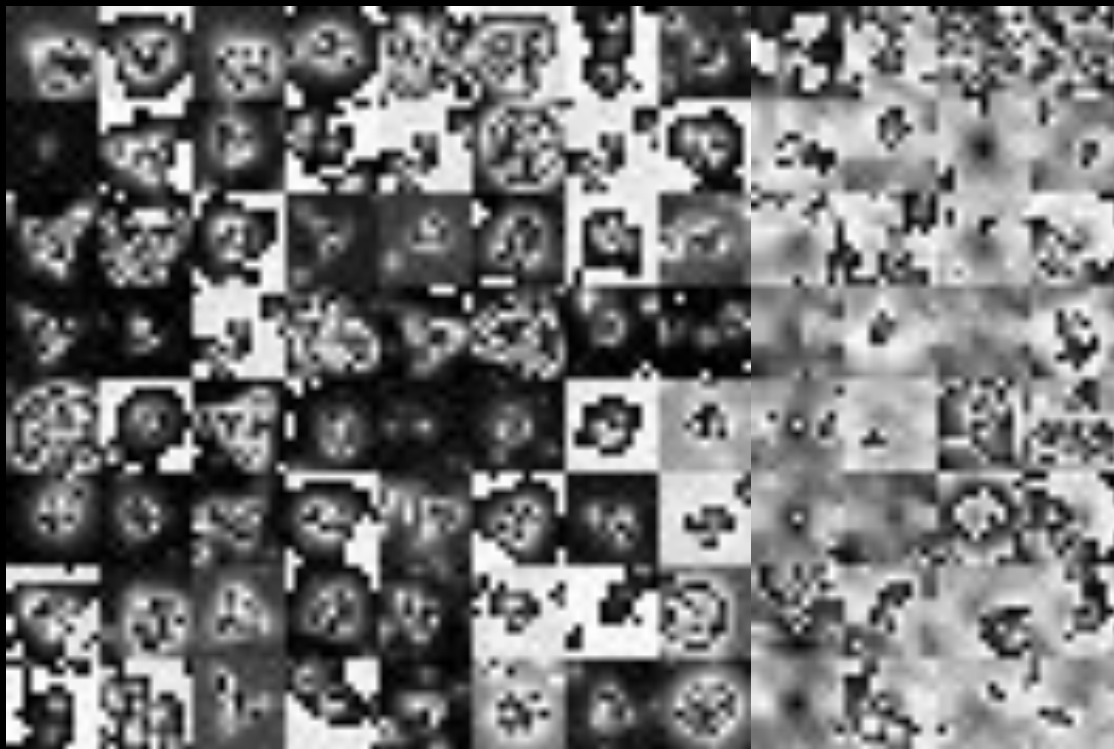Stellar variability

Binary

Planet

FDL

# From pixels...

64 000 month-long videos

Lots of noise!

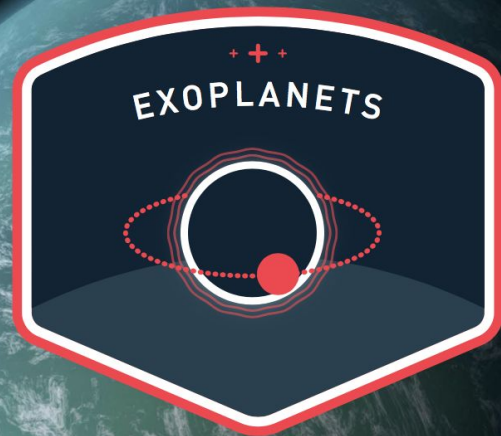Promising approach which could by-pass TESS pipeline…

But needs more work.

# Summary

- Classify TESS planets **faster** & **more precisely** than previous approaches.

- Innovative new avenues for planet hunting direct **from light curves & pixels**
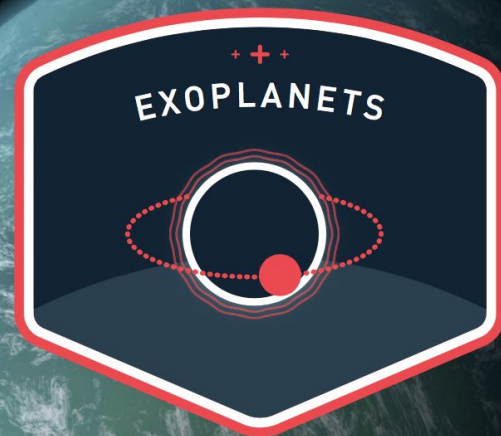
**Megan Ansdell**, Yani Ioannou, **Hugh Osborn**, Michele Sasdelli + Jeff Smith, Jon Jenkins, Doug Caldwell
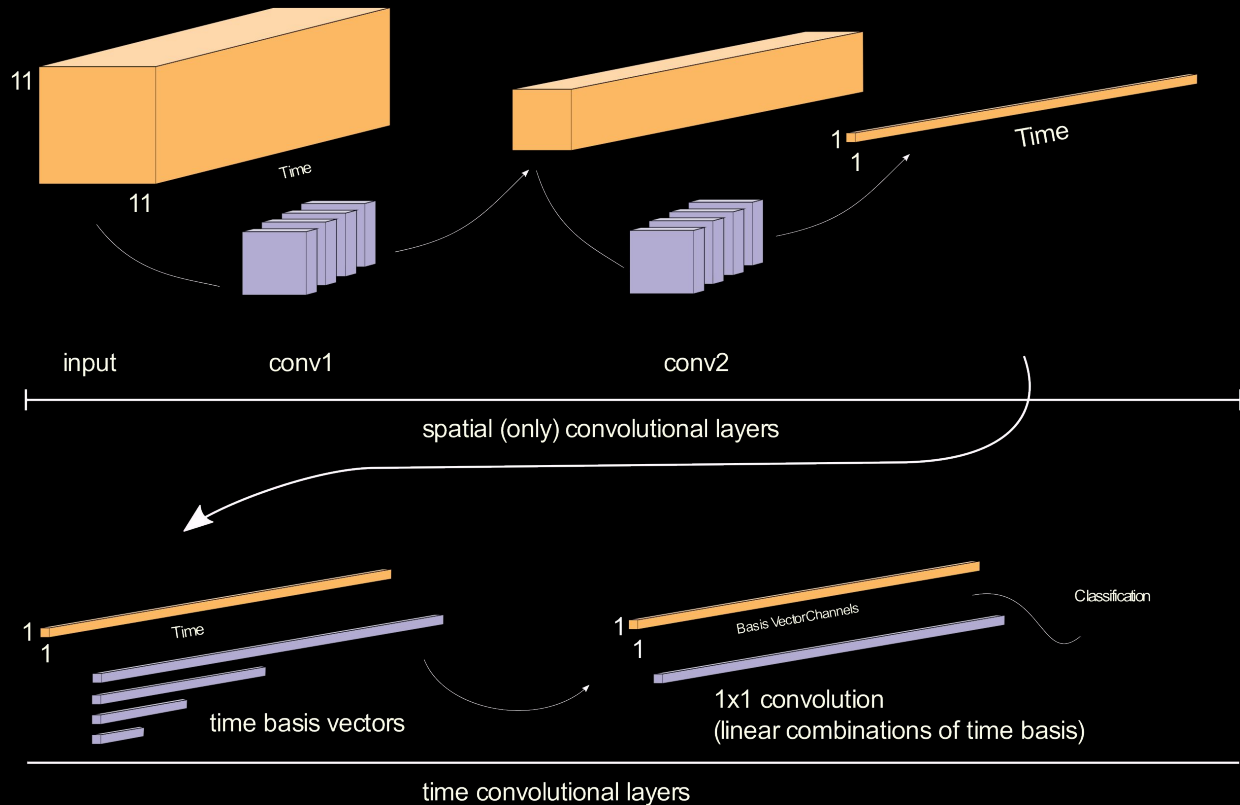
*Thanks*

**Megan Ansdell**, Yani Ioannou, **Hugh Osborn**, Michele Sasdelli
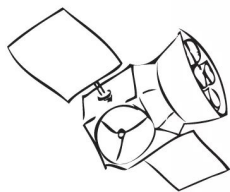+ Jeff Smith, Jon Jenkins, Doug Caldwell

# From pixels...

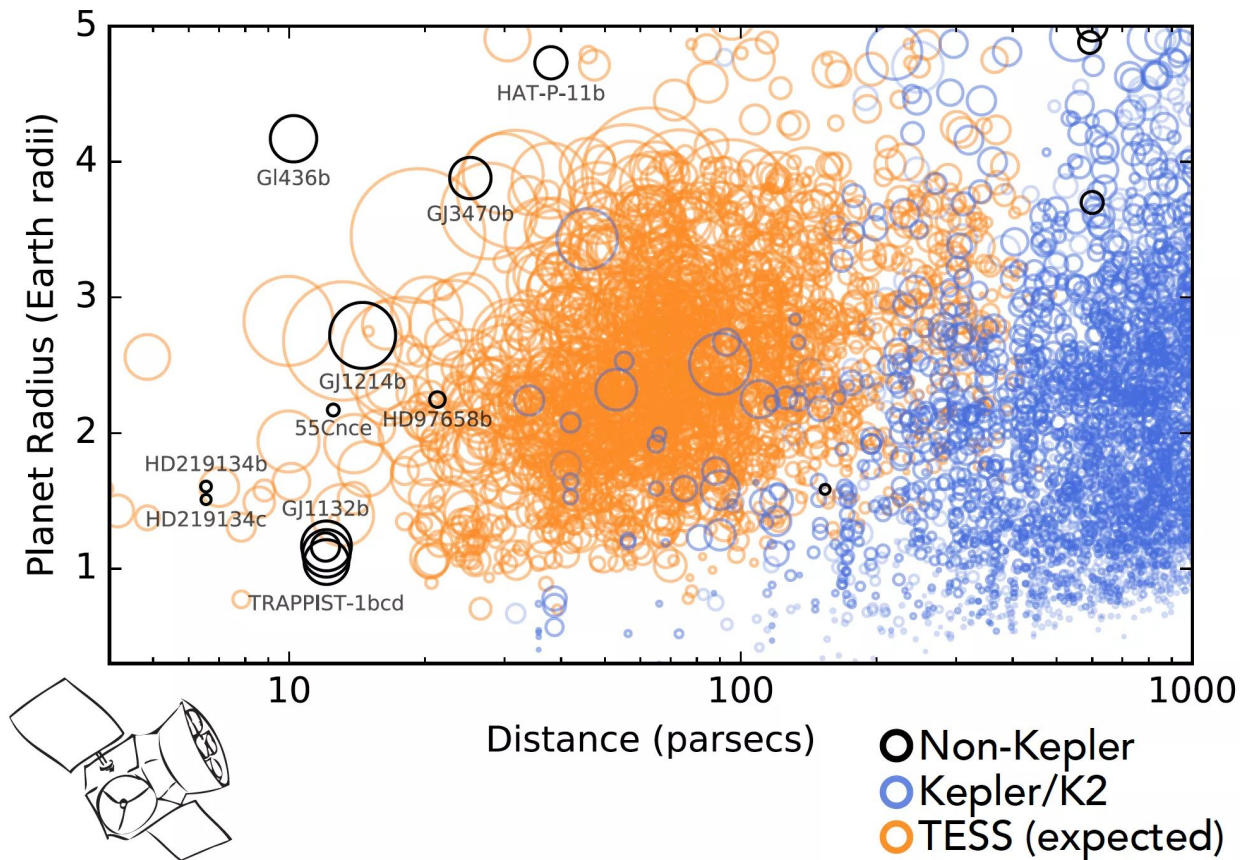64 000 videos of stars
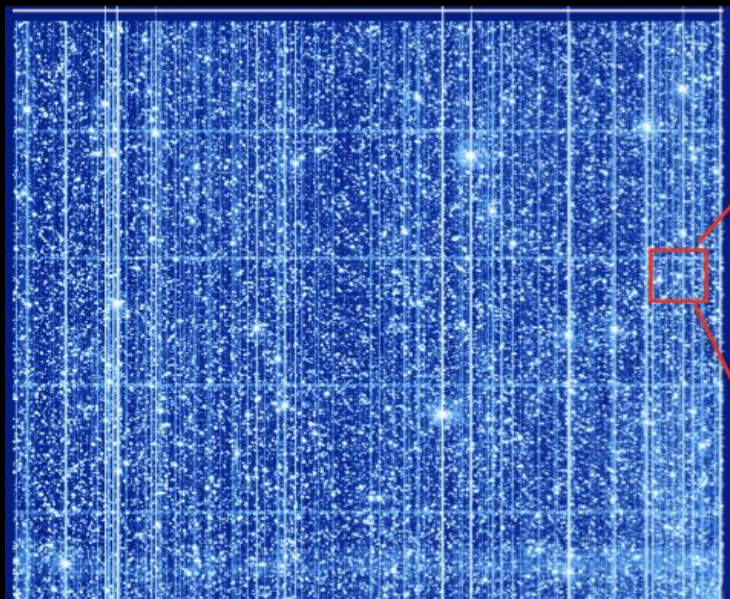
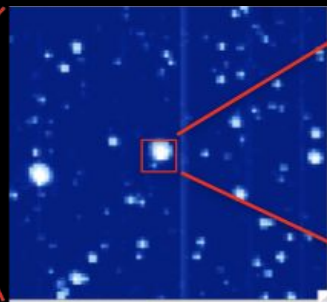Promising approach which could by-pass TESS pipeline…

But needs more work



11

11

Time

1
1

Time

input          conv1                          conv2

spatial (only) convolutional layers

1
1

Time

time basis vectors

1
1

Basis VectorChannels

Classification

1x1 convolution
(linear combinations of time basis)

time convolutional layers

FDL

intel AI    SR SPACE RESOURCES.LU    XPRIZE    Google Cloud    nVIDIA    LOCKHEED MARTIN    kx    IBM    KBRwyle We Deliver

# The data: TESS mission

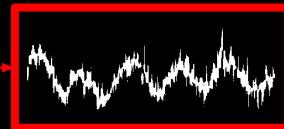# TESS

Exoplanets
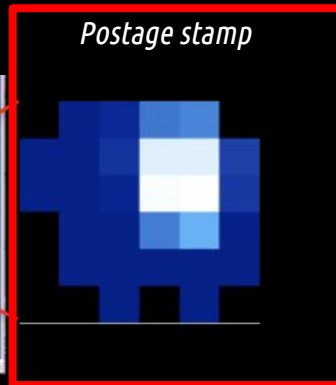around nearby
stars

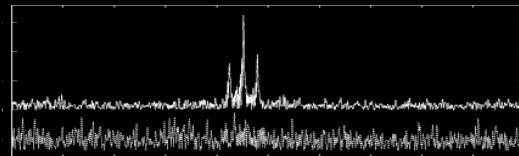Full-frame images

Zoom-in

Postage stamp

# Incremental: Classify planet candidates detected by the pipeline

- Train a Neural Network on candidate planets (TCEs) detected by the pipeline Transiting Planet Search

- Use heavily pre-processed domain data generated by the pipeline.

*In Shallue et al 2017:*

*Full lightcurve*

*Zoomed-in lightcurve*

*New Input Datasets:*

*Centroids (x & y)*

*Stellar Properties*

*Data augmentation*

# Incremental: Classify planet candidates detected by the pipeline

Improvements on Shallue, 2017:

- More input datasets

- Multiple false-positive classes

- Use data balancing

- Augment light curves to increase training examples.

- Lighter NN

- 

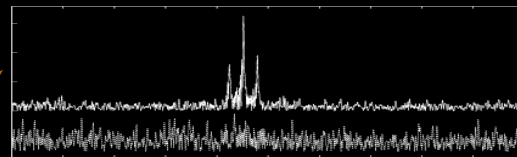*In Shallue et al 2017:*

*Full lightcurve*

*Zoomed-in lightcurve*
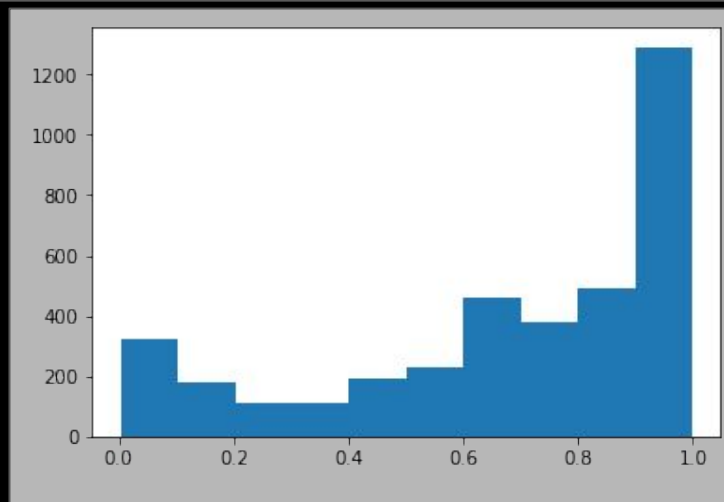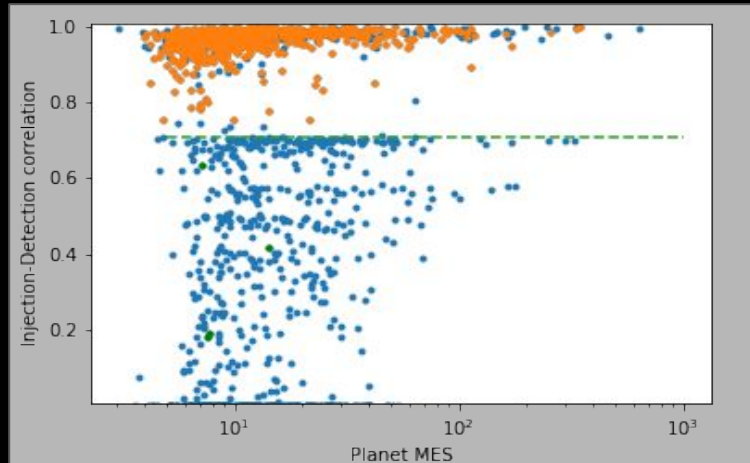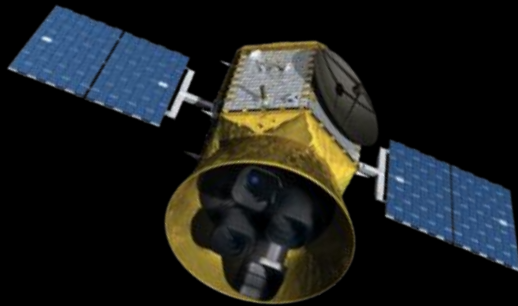
*New Input Datasets:*

*Centroids (x & y)*

*Stellar Properties*

*Frequency-space data*

# Ongoing work:

- Labels not necessarily intuitive.
- Modified "correlation" metric between injections and detections to include period multiples.
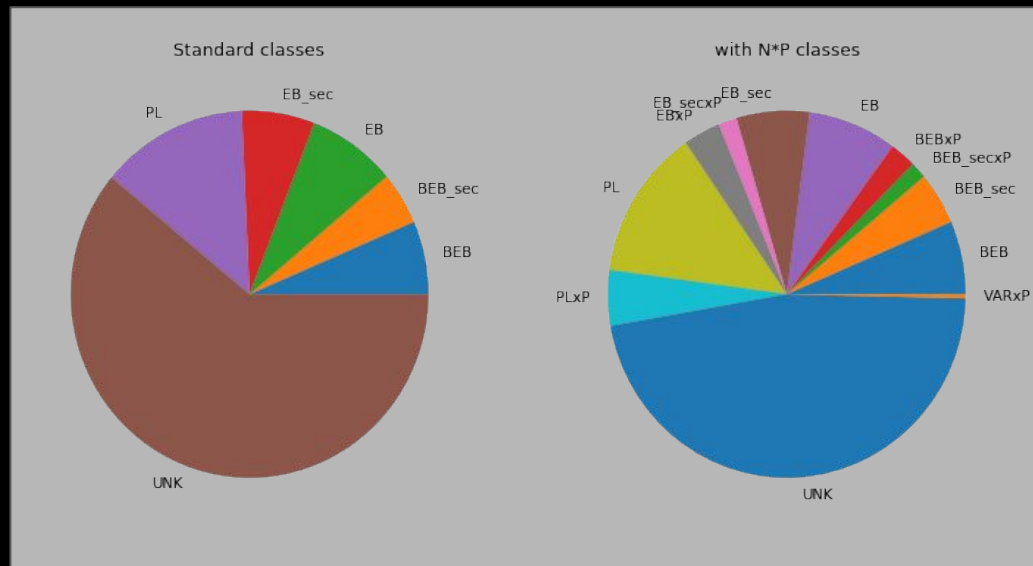- Still issues with high-SNR planets not being detected
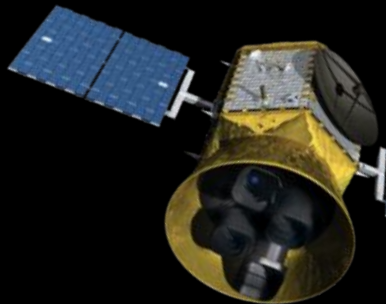
# Ongoing work:

- Labels not necessarily intuitive.
- Modified "correlation" metric between injections and detections to include period multiples.
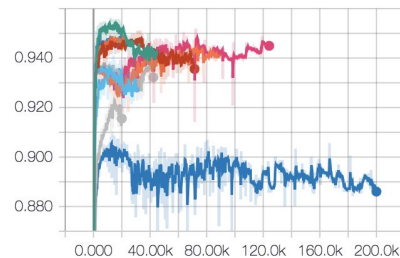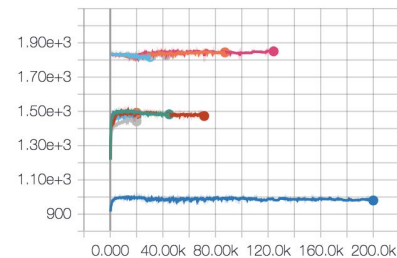- Still issues with high-SNR planets not being detected

# Ongoing work:

- Converted Shallue code from Kepler to TESS
- Improvement with centroids
- Also improved with modified smoothing (spline) techniques.
- Models being trained
- Stellar parameters and frequency space data to be tested and added

# Next steps...

Week 4:

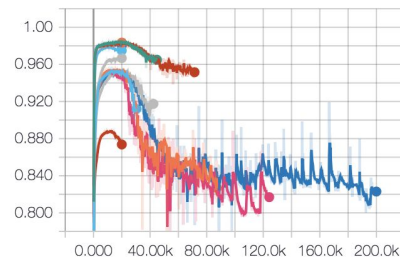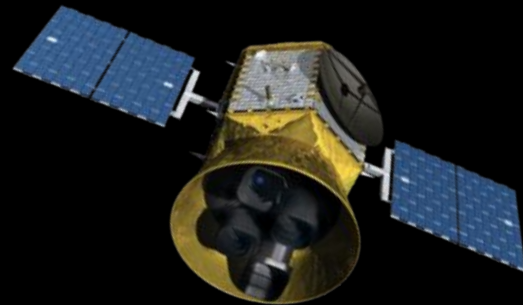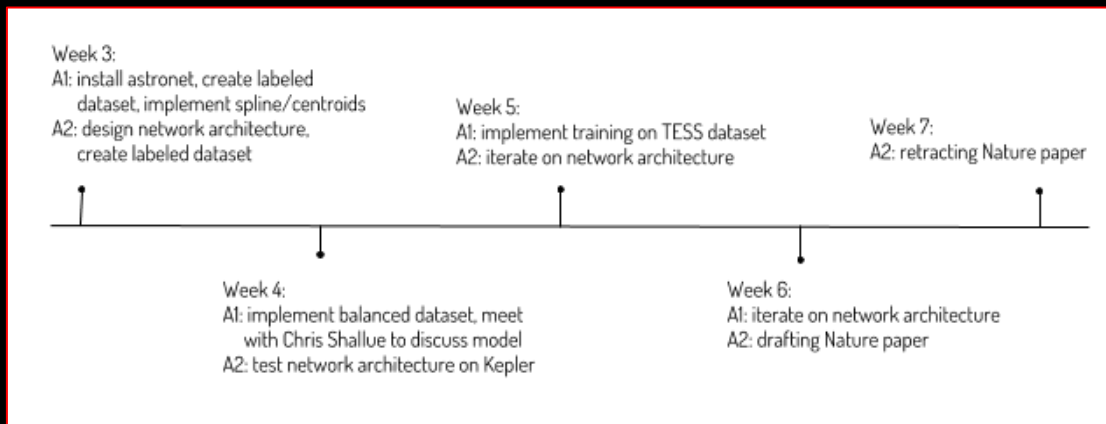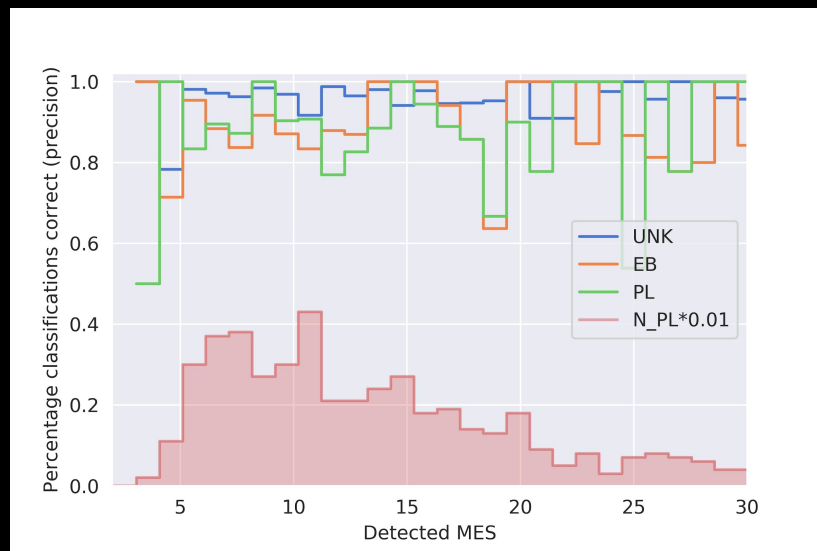- *Avenue 1*: Implement balanced labels
- *Avenue 2*: Test design architecture with Kepler

Week 5:

- *Avenue 1*: Implement training based on TESS data
- *Avenue 2*: Iterate on model/architecture

Week 3:
A1: install astronet, create labeled
    dataset, implement spline/centroids
A2: design network architecture,
    create labeled dataset

Week 5:
A1: implement training on TESS dataset
A2: iterate on network architecture

Week 7:
A2: retracting Nature paper

Week 4:
A1: implement balanced dataset, meet
    with Chris Shallue to discuss model
A2: test network architecture on Kepler

Week 6:
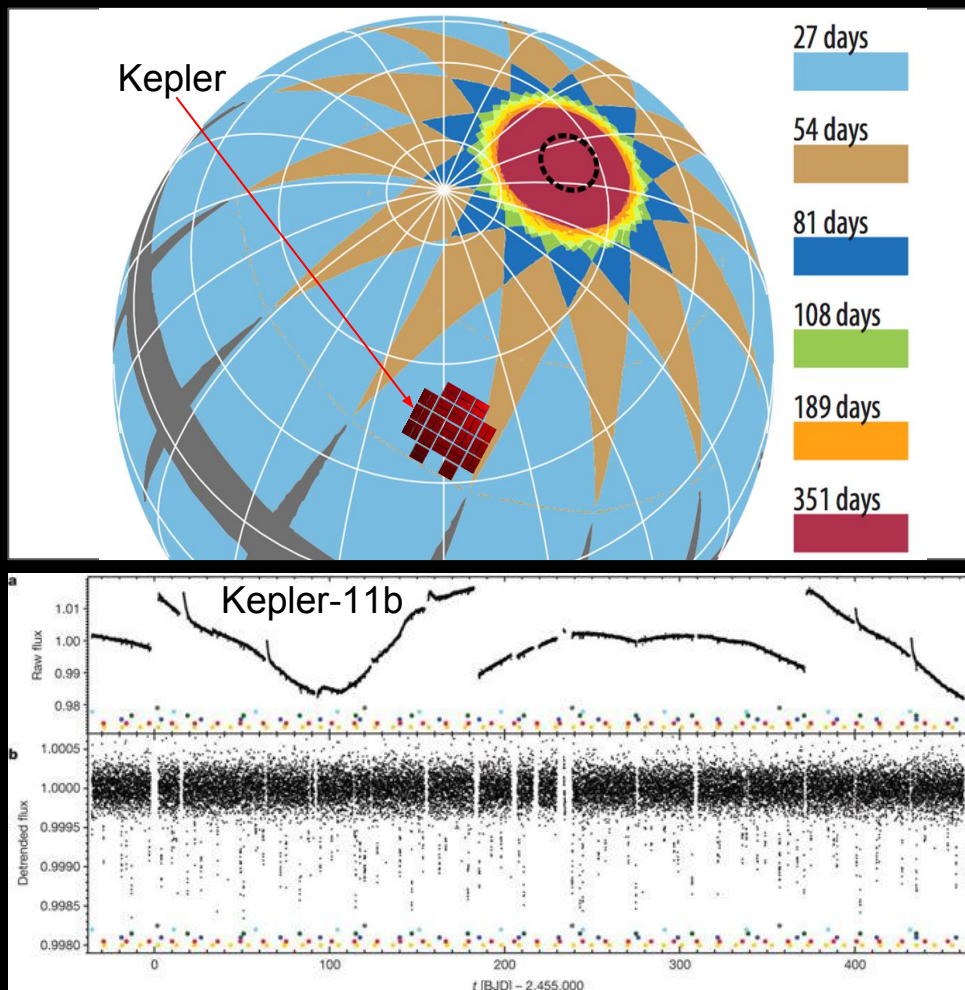A1: iterate on network architecture
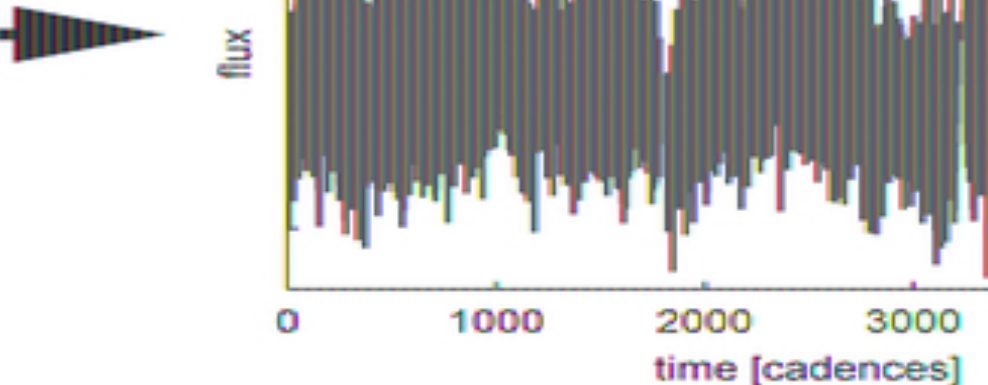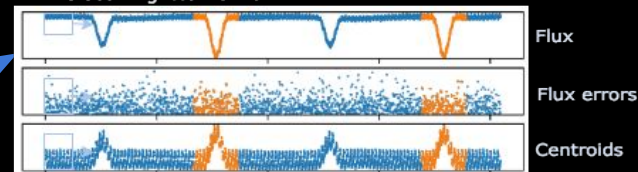A2: drafting Nature paper

# Results on TESS

# Kepler Data

Real data but noisy labels

- 150 000 lightcurves

- Lightcurves 4 years in duration

- ~4 000 planets (candidates & confirmed)

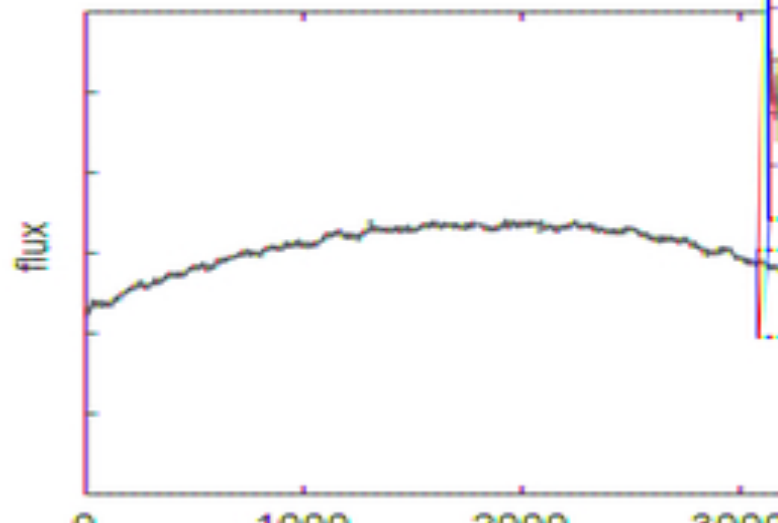- Can augment data to TESS-like 27-day campaigns (up to 7.8 million 'targets')



Kepler

27 days
54 days
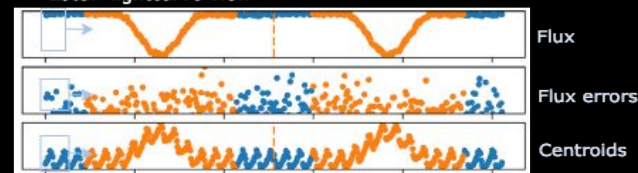81 days
108 days
189 days
351 days



Kepler-11b

flux

time [cadences]

C

fast oscillation

flux

"Global" lightcurve view

Flux

Flux errors

Centroids

"Local" lightcurve view

Flux

Flux errors

Centroids
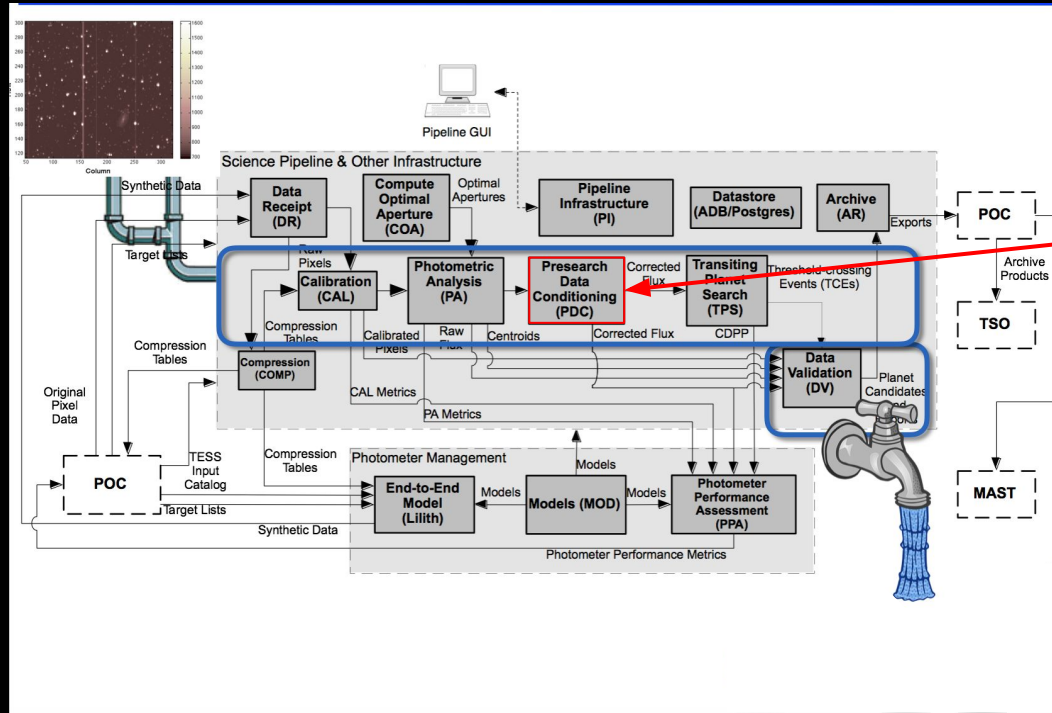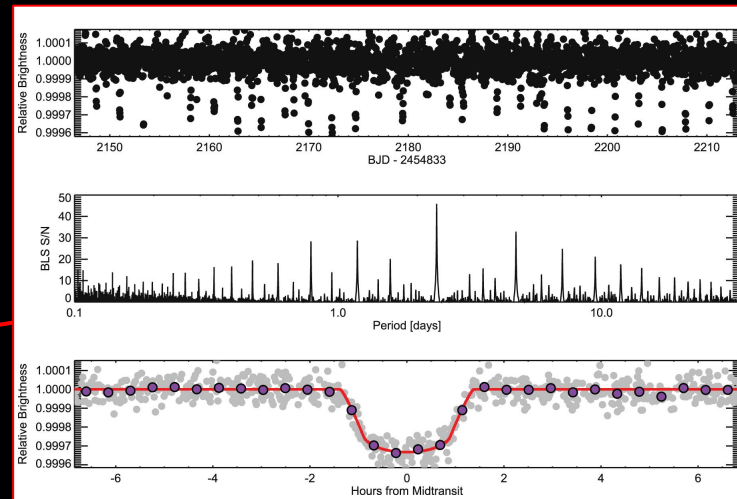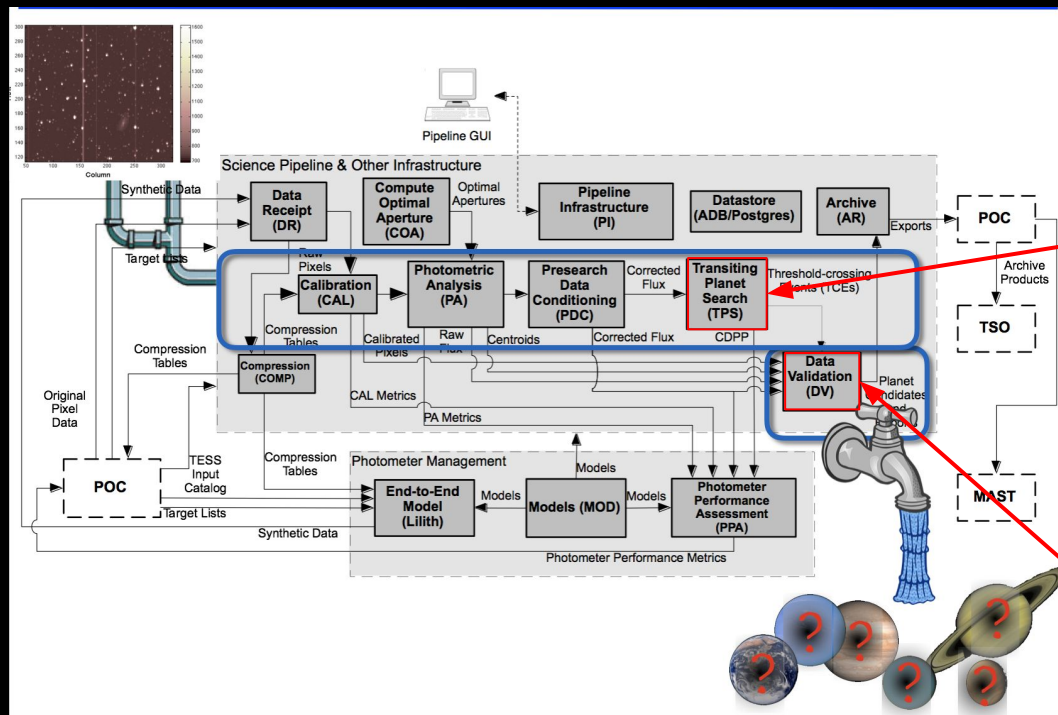
# The data: real light curves

# The data: detrended light curves



- Remove trends present on all stars to leave only astrophysical signal

# The data: Transiting Planet Search



- Search for transiting planets in frequency domain.
- Strong candidates are analysed using statistical tests (DV)

# The data: planet candidates



Planet candidates (TCEs) are analysed to determine which are planets.

Often performed:
- By human vetters
- Using classical statistical techniques
- Using Machine Learning

...Could we start from this point?

# Incremental: Classify planet candidates detected by the pipeline



Classify candidates
detected by TESS pipeline

*Planet, EB, BEB?*

# Innovative approach: Classifying directly from pixel data



Target Pixel File

Cut out the
pipeline entirely

*Planet, EB, BEB?*

# Innovative approach: Classifying directly from pixel data
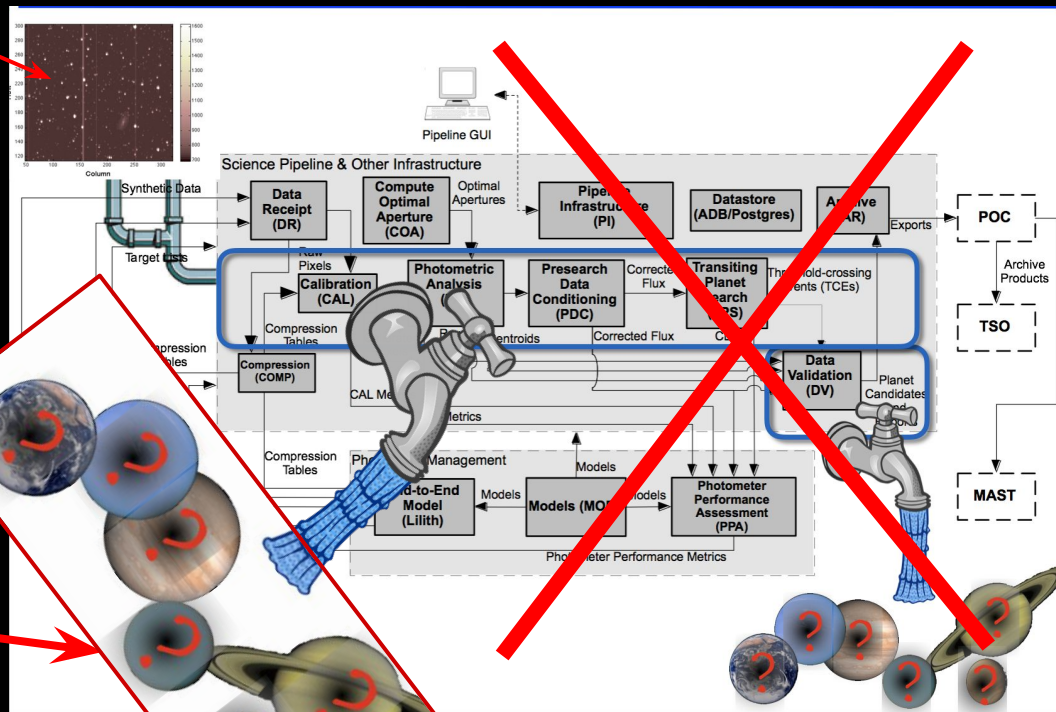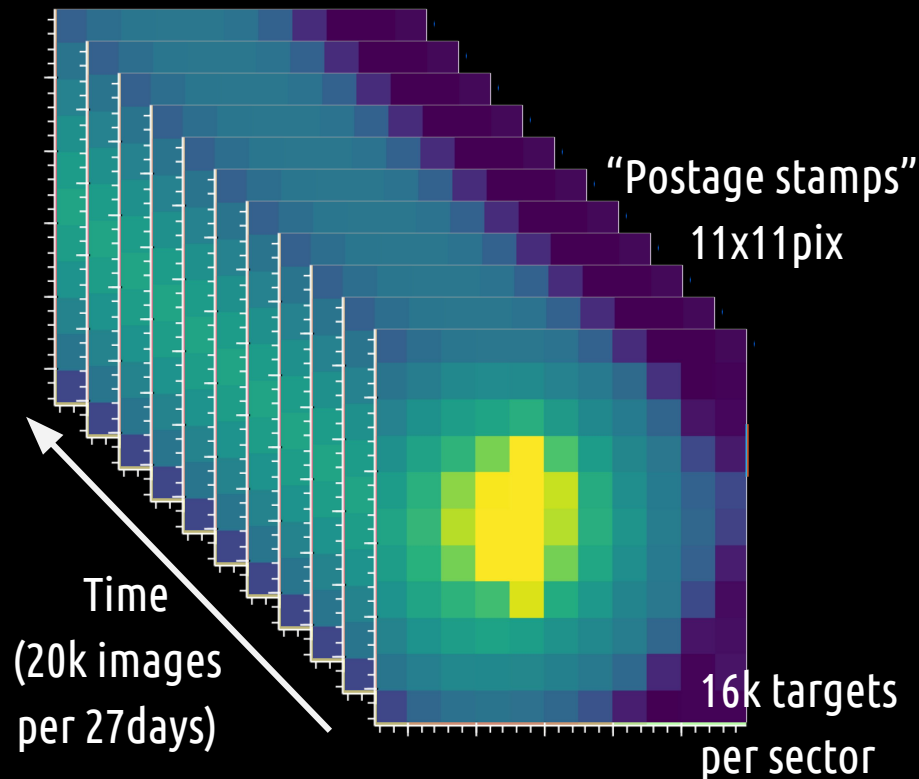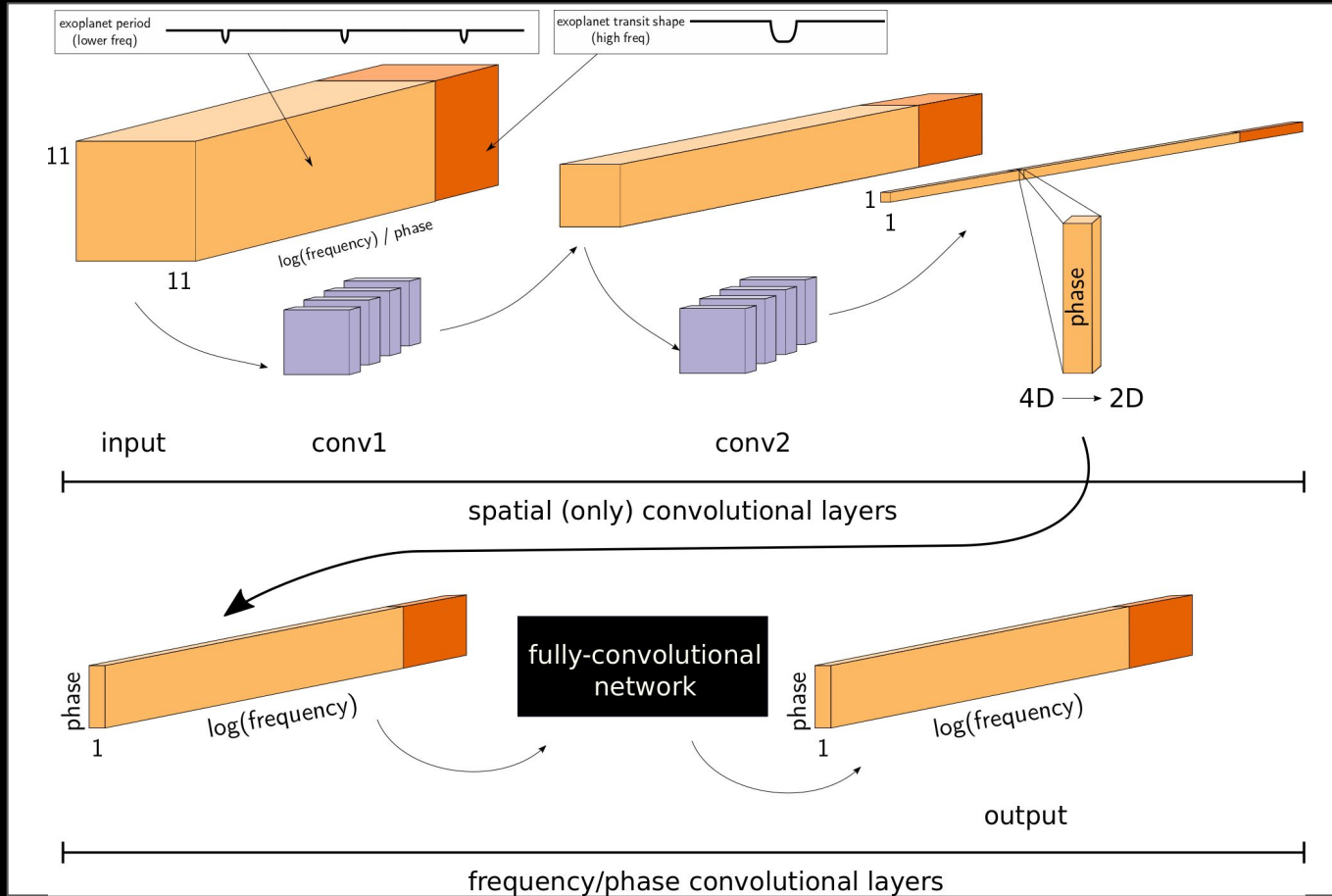
- Train a Convolutional Neural Network (CNN) directly on pixels

- Transform image to more appropriate representation, i.e. frequency domain for periodic data

- Incorporate our domain knowledge into the network architecture and learning algorithm, rather than the data



"Postage stamps" 11x11pix

Time (20k images per 27days)

16k targets per sector

FDL

SETI INSTITUTE · intel AI · SR SPACE RESOURCES.LU · XPRIZE · Google Cloud · nVIDIA · LOCKHEED MARTIN · kx · IBM · KBRwyle We Deliver

exoplanet period
(lower freq)

exoplanet transit shape
(high freq)

11

11

log(frequency) / phase

1
1

phase

4D → 2D

input          conv1                    conv2

spatial (only) convolutional layers

phase

1

log(frequency)

fully-convolutional
network

phase

1

log(frequency)

output
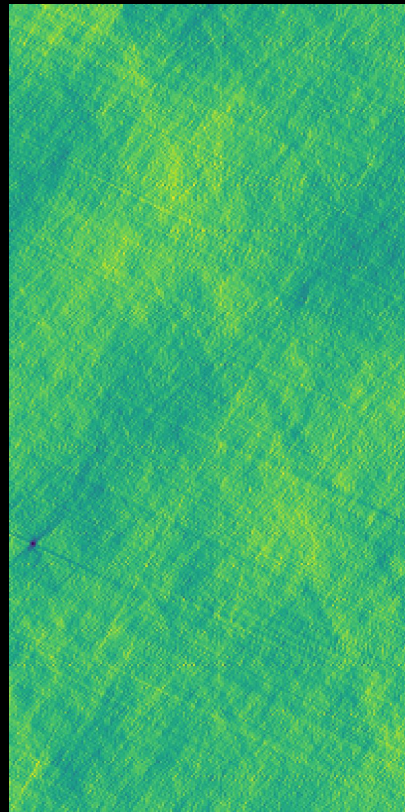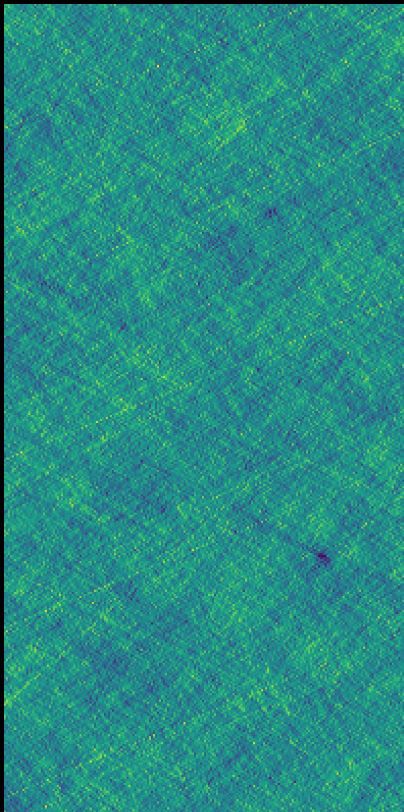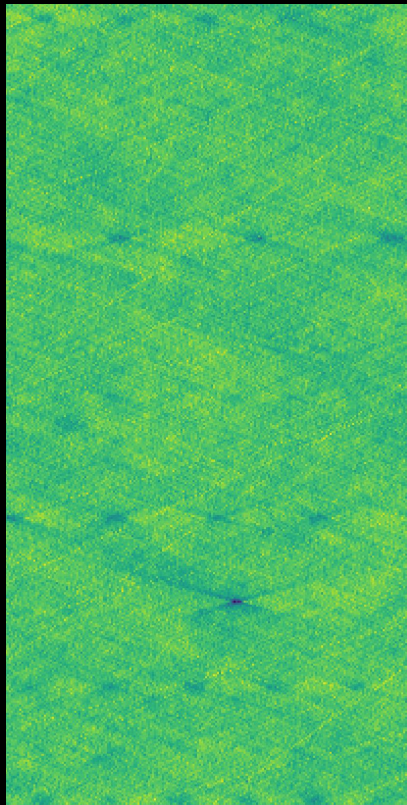
frequency/phase convolutional layers

# Next steps...

Week 3:

- *Avenue 1*: Re-produce Shallue & Vanderburg 2017 results
- *Avenue 2*: Setup data infrastructure/get basic CNN training

Week 4:

- *Avenue 1*: Incorporate new domain knowledge
- *Avenue 2*: Iterate on model/implement data balancing methods

# Planets

# No planets