

# AstroCompress:

A benchmark dataset for multi-purpose  
compression of astronomical data

**Tuan Truong\*, Rithwik Sudharsan\***

Yibo Yang, Peter Ma, Ruihan Yang

Stephan Mandt, Joshua Bloom

\*Equal contribution first authors. Direct correspondence to:

<tuannt2@uci.edu>, <rithwik@berkeley.edu>



# 01: The Problem

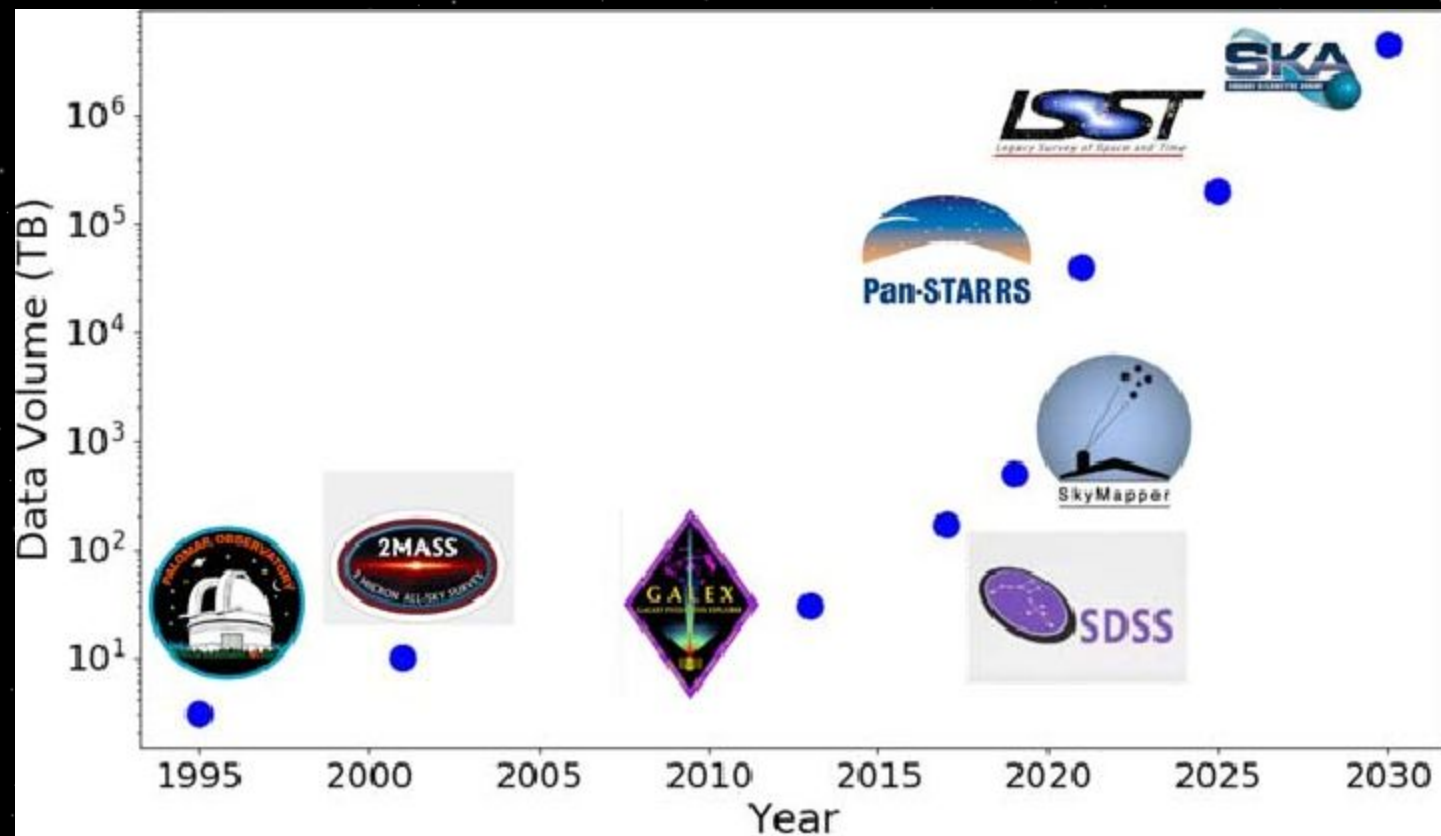
The state-of-the-art James Webb Space Telescope (JWST) has:

- ~65 GB storage capacity
- ~60 GB is downlinked per day at 28 Mbps
- Contact with Earth exists only 8 hours per day → lots of downtime

JWST is at L2, but even satellites in LEO (like CuRIOS-ED) struggle with bandwidth.

There certainly is a desire to collect more data: at the extreme, ground-based telescopes like the Square Kilometer Array will collect ~1 exabyte per year / 3 Tbps, which is 40,000x JWST.

- Transporting this data even on the ground is a challenge, so compression is needed!



## SDP Challenges:



Input

- 400 Gbyte/s INGEST

Distribute

- 400 million tasks in the graph
- Half an Exaflop/s total peak

Generate  
(and destroy)

- 1.3 ZettaBytes intermediate data products

Preserve and  
ship

- 1 PetaByte per day of Science Data Products
- Each cube up to a few PBytes

# 02: Shannon's Theorem

Min compression bitrate = entropy of X.

Density estimation = lossless compression.

## Entropy

### Entropy, Cross-Entropy, and KL-Divergence

- The Kullback-Leibler or “KL” Divergence is defined by

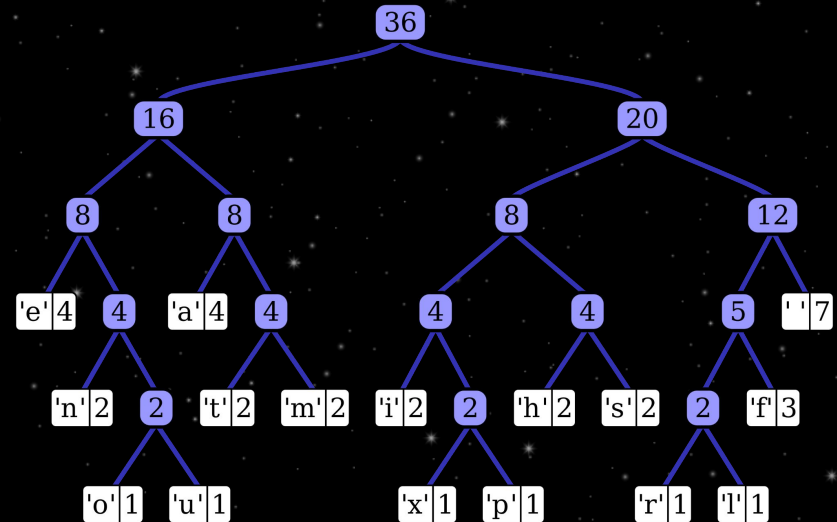
$$KL(p||q) = \mathbb{E}_p \log \left( \frac{p(X)}{q(X)} \right).$$

- $KL(p||q)$ : #(**extra bits**) needed if we code with  $q(x)$  instead of  $p(x)$ .
- The **cross entropy** for  $p(x)$  and  $q(x)$  is defined as

$$H(p, q) = -\mathbb{E}_p \log q(X).$$

- $H(p, q)$ : #(**bits**) needed to code  $X \sim p(x)$  using  $q(x)$ .
- Summary:

$$H(p, q) = H(p) + KL(p||q).$$



# 03: Existing Codecs

## RICE (used in astronomy today)

- Lightning fast, almost no spatial information used, entropy coding assuming Geometric distribution

## HCOMPRESS, JPEG-2000

- Wavelet transforms + entropy coding
- Sometimes use other things like run-length encoding and LZ77

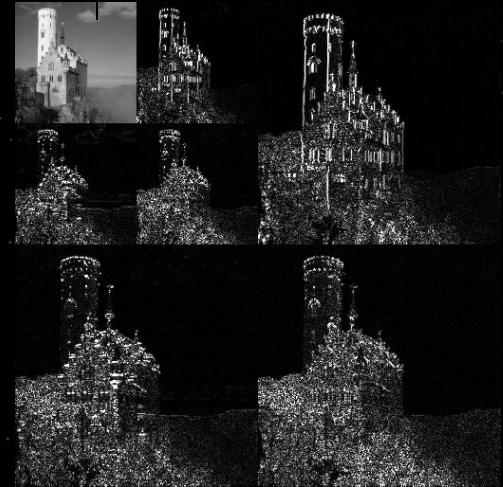
## JPEG-LS

- Autoregressive prediction w/ context modeling + entropy coding of residual

## JPEG-XL

- All of the above on steroids

	C	B	D	
	A	X		



# 04: Our Dataset

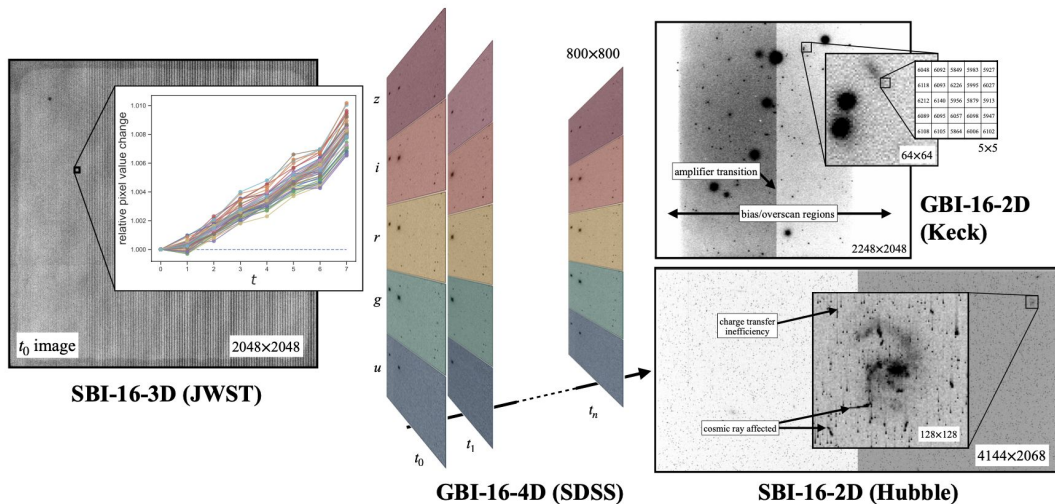


Figure 1: Depiction of salient features in the AstroCompress corpus using representative images from each dataset. Inset to the JWST  $t_0$  (first) image are the value changes in time for a small sample of pixels. In SDSS there are 5 filtered images per observation epoch, up to a variable number  $n$  observations in the same portion of the sky. The inset of Hubble zooms in on a spiral galaxy, showing cosmic ray hits (black) and charge transfer inefficiency, causing vertical flux smearing. The actual pixel values in Keck are shown for a zoomed-in  $5 \times 5$  pix<sup>2</sup> region.

# 05: ML Methods: IDF

- Normalizing flows
- Integer operations only, to avoid lossy quantization at the end
- Factored logistic mixture distribution enforced on  $z$
- During decoding, we only have to decode the latent  $z$  based on the logistic mixture
- Incentive to decorrelate  $z_i$  is implicit in the “factored” distribution.  
That’s why we don’t just compress  $x$  instead of  $z$
- Jacobian = 1
- Exact density  $p(z)=p(x)$  is known

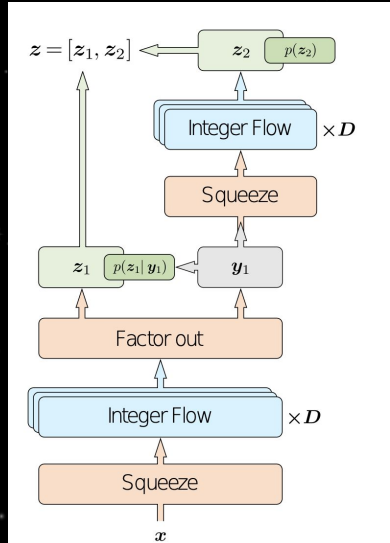
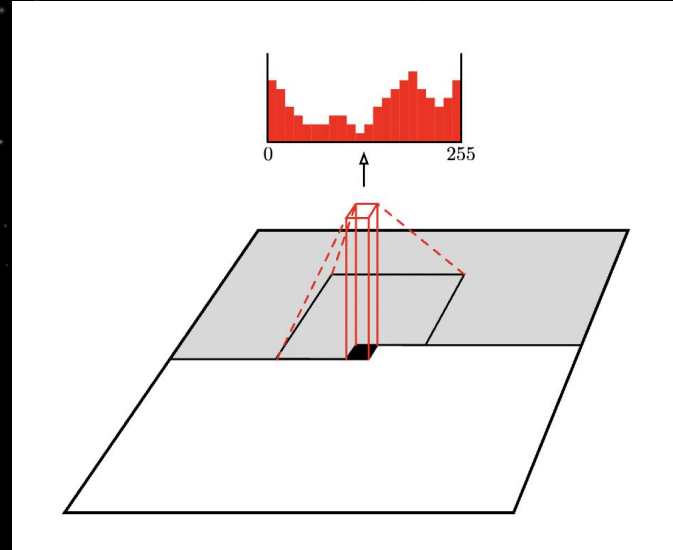


Figure 4: Example of a 2-level flow architecture. The squeeze layer reduces the spatial dimensions by two, and increases the number of channels by four. A single integer flow layer consists of a channel permutation and an integer discrete coupling layer. Each level consists of  $D$  flow layers.



# 05: ML Methods: PixelCNN++

- Autoregressive
- Predict distribution over possible values  
    Again, logistic mixture (discretized)
- Simple, similar to JPEG-LS
- Parallelizable encoding, slow decoding
- Exact density  $p(x)$  is known



# 05: ML Methods: L3C

- Basically an autoregressive model that models **features at different scales** autoregressively rather than every single pixel
- Thus, it's super fast
- Parallels to U-Net:  $D(k)$  passes skip connections to  $D(k-1)$
- It's cool and creative but not that effective

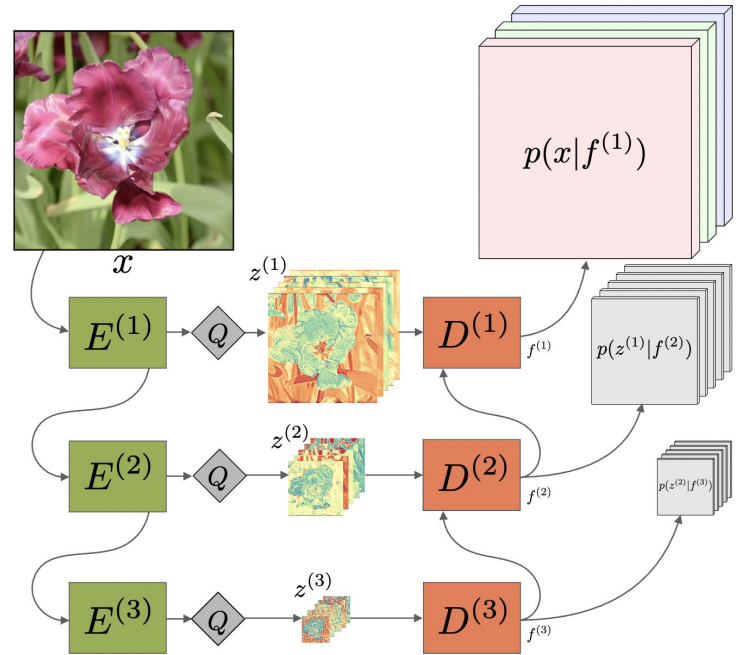
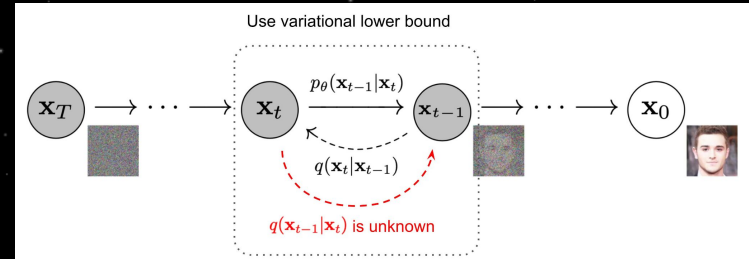


Figure 1: Overview of the architecture of L3C. The feature extractors  $E^{(s)}$  compute quantized (by  $Q$ ) auxiliary hierarchical feature representation  $z^{(1)}, \dots, z^{(S)}$  whose joint distribution with the image  $x$ ,  $p(x, z^{(1)}, \dots, z^{(S)})$ , is modeled using the non-autoregressive predictors  $D^{(s)}$ . The features  $f^{(s)}$  summarize the information up to scale  $s$  and are used to predict  $p$  for the next scale.

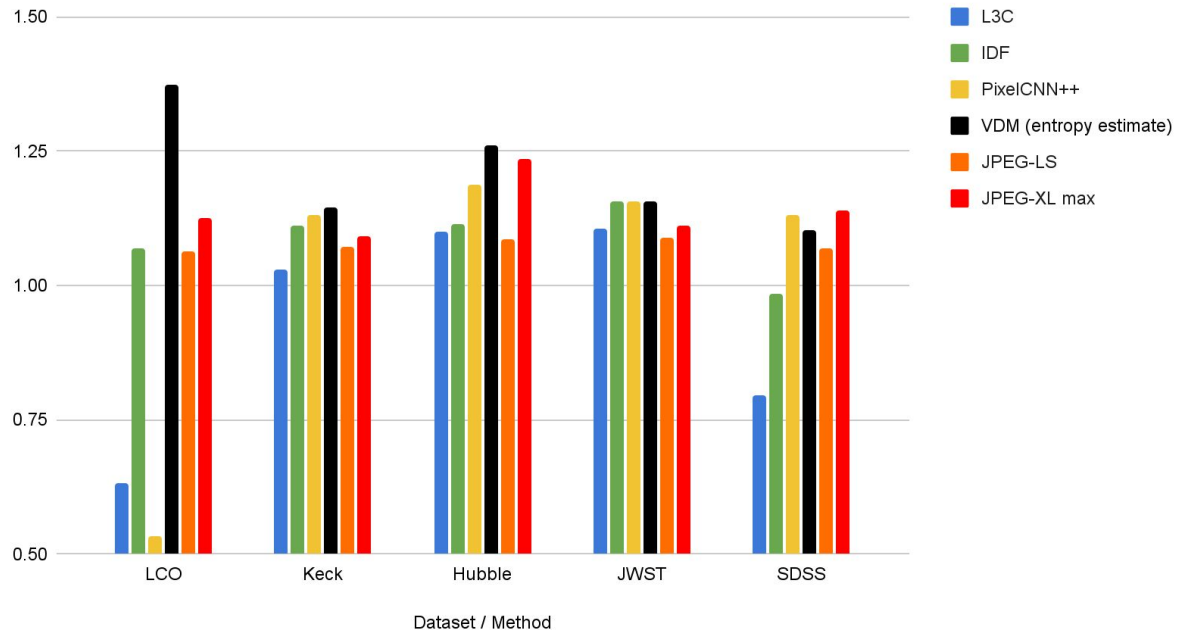
# 05: ML Methods: VDM

- Imagine 1 VAE... a VDM is ~a bunch of stacked VAE's where each one is predicting the slightly denoised image instead of reconstructing the whole image
- Issue: latent is always a sample of a random distribution. Need "bits-back" coding to actually do compression
- Transmit the noisiest latent  $z_1$  with prior  $p(z_1)$
- Encode residuals at each stage of VDM with  $p(\text{current } z \mid \text{previous } z)$
- It's like a diffusion model where we always have the exact correct image at each stage
- Extremely slow (but optimizations exist)



# 06: Results

Compression Ratios, relative to RICE (current astronomy paradigm)



# 06: Results

Codec	SDSS-2D (800x800)	Hubble (4144x2068)
IDF	$0.42 \pm 0.01$	$6.03 \pm 0.24$
L3C	$5.18 \pm 1.04$	$73.04 \pm 2.36$
PixelCNN++	$1.48 \pm 1.05$	$20.49 \pm 0.18$
JPEG-XL max	$3.14 \pm 0.14$	$87.76 \pm 13.30$
JPEG-XL default	$0.06 \pm 0.002$	$0.91 \pm 0.07$
JPEG-LS	$0.02 \pm 0.0002$	$0.316 \pm 0.04$
JPEG-2000	$0.09 \pm 0.003$	$1.76 \pm 0.11$
RICE	$0.008 \pm 0.0002$	$0.12 \pm 0.02$
VDM	$2301 \pm 229.2$	$33072 \pm 19.8$

Table 3: Compression (encoding) runtime (in seconds/image) on the SDSS-2D and Hubble datasets. For neural methods, we measure the time for evaluating the likelihood under the model without entropy coding.

Training Set	Evaluation Set				
	LCO	Keck	Hubble	JWST-2D	SDSS-2D
LCO	<b>2.83</b>	1.01	1.09	0.84	2.31
Keck	2.70	<b>2.05</b>	2.20	1.19	<u>3.02</u>
Hubble	0.67	0.94	<u>2.94</u>	1.22	0.69
JWST-2D	1.46	1.45	1.47	<b>1.44</b>	1.50
SDSS-2D	2.27	1.24	1.75	1.02	2.91
All data	<u>2.82</u>	<u>1.87</u>	<b>2.98</b>	<u>1.38</u>	<b>3.18</b>

Table 2: IDF generalized performance across single-frame datasets. Rows indicate train set; columns indicate test set. Bold indicates best in test set; underline indicates second-best.

# 07: Analysis

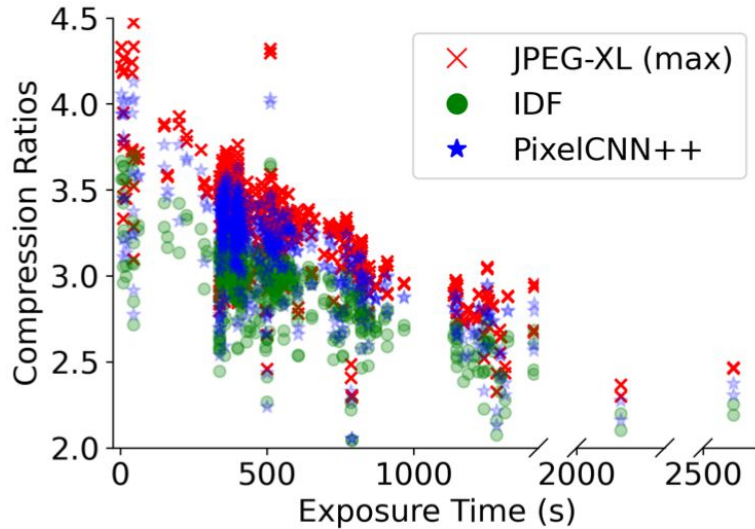


Figure 2: Hubble exposure times plotted against compression ratios using various algorithms. Longer exposure times tend to induce more incompressible noise and, hence, reduce compression ratios.

# 07: Analysis

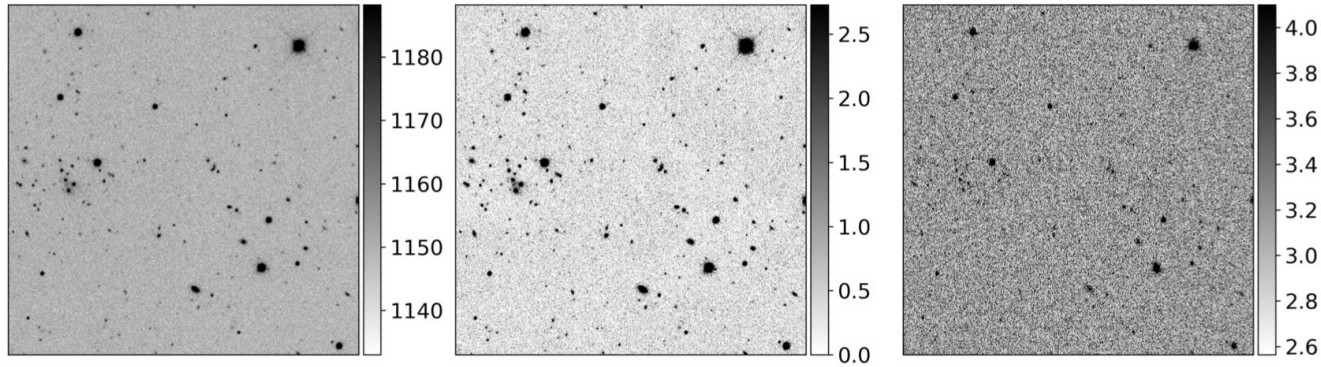
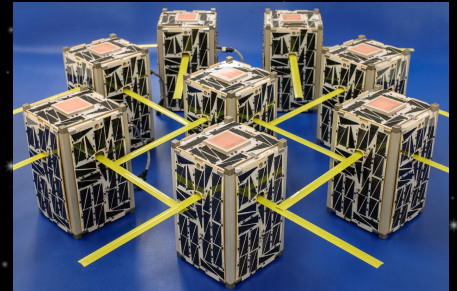


Figure 3: From left to right, an example SDSS-2D image: raw image, SNR heatmap, and PixelCNN++ bitrate heatmap. Colors are z-score normalized for visualization; colorbars indicate true values.

# 08: Future Directions

- **Near-lossless compression**
- **CuRIOS satellite**
- **Autoregressive models' encode-decode time tradeoff**
- **More work on 3D / 4D data cube compression**



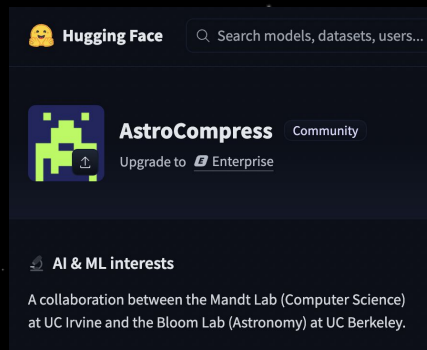


# Published! (to an ML Audience)

Under review as a conference paper at ICLR 2025

000 **ASTROCOMPRESS:**  
001  
002 **A BENCHMARK DATASET FOR MULTI-PURPOSE**  
003 **COMPRESSION OF ASTRONOMICAL IMAGERY**  
004

005  
006 **Anonymous authors**  
007 Paper under double-blind review  
008



External sources can be found by clicking any image in this slides file, which can be found at: <https://github.com/tuatruog/AstroCompress>