

SEARCH

A graphic element for the 'SEARCH' logo. It consists of a square tilted at an angle, filled with a golden, turbulent, solar-like texture. This square is enclosed within three concentric, slightly offset square borders, also tilted at the same angle. The innermost border is a thin gold line, the middle is a thicker dark blue line, and the outermost is a thin gold line.

SDO Exploration And Research Community for Heliophysics

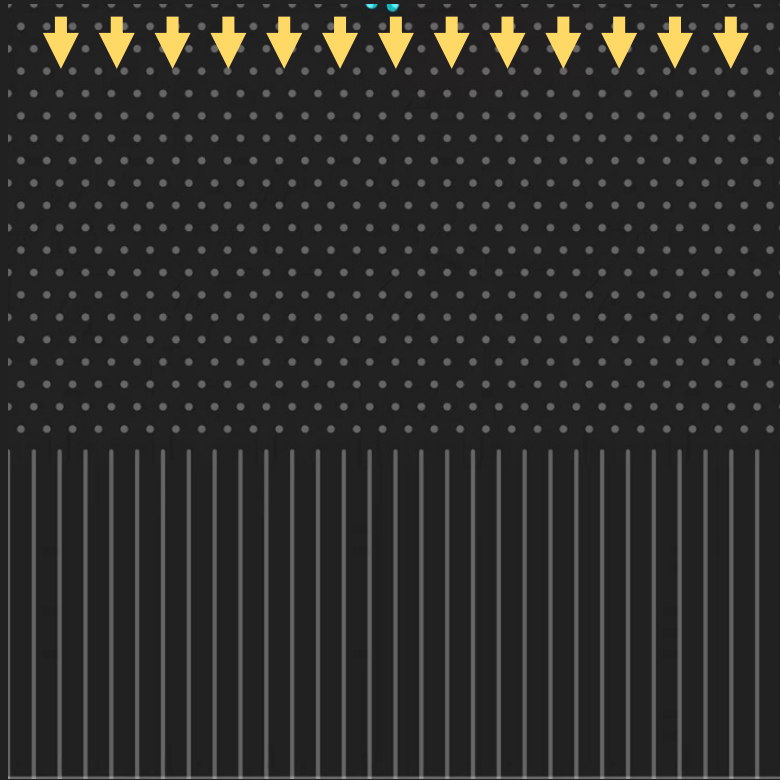
An experiment on using applied AI research to educate a
diverse workforce

Subhamoy Chatterjee, Nadia Ahmed, & Andrés Muñoz-Jaramillo
andres.munoz@swri.org

Funded by NASA-HITS



Why SEARCH?



Our lives are very random. They are **NOT** the deterministic outcome of the choices we make.

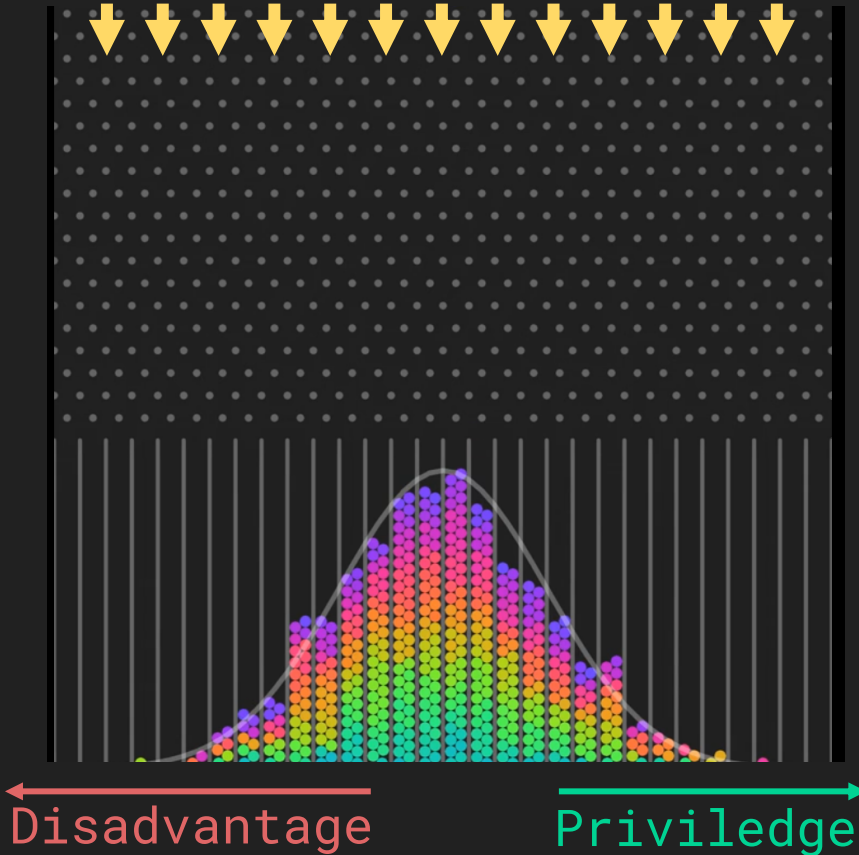
Our birth has an outsized influence on the rest of our lives. We do **NOT** start with the same circumstances and opportunities.

Disadvantage

Privilege



Why SEARCH?

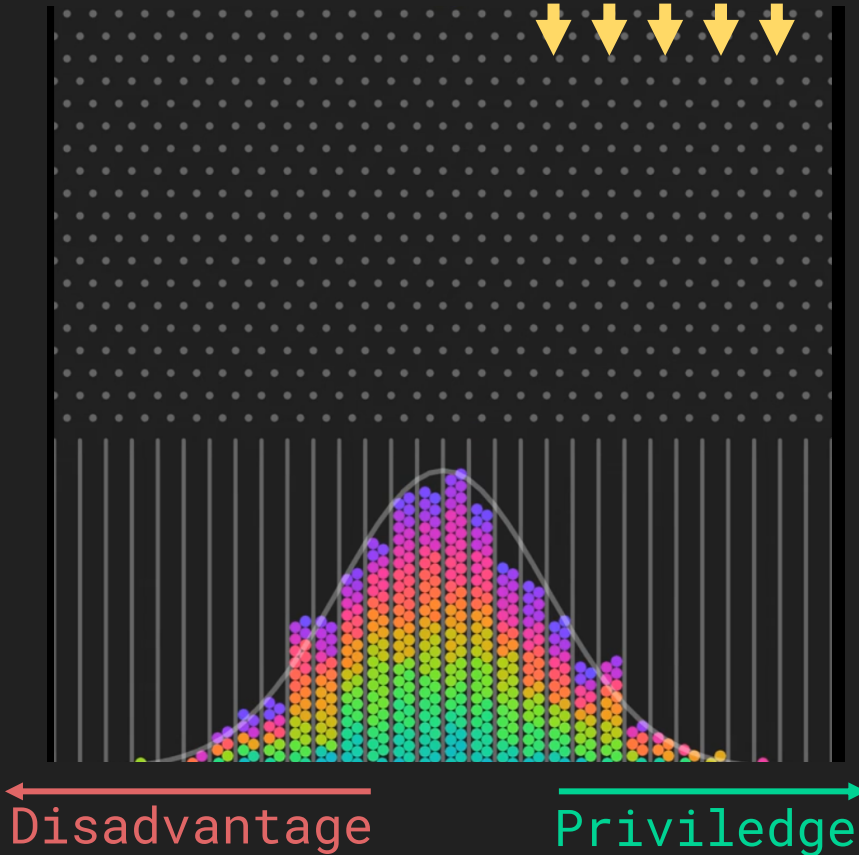


Education is one of the main ways we have to help disadvantaged populations.

Most educational opportunities are gated, which tends to help those that are ahead get farther ahead.



Why SEARCH?

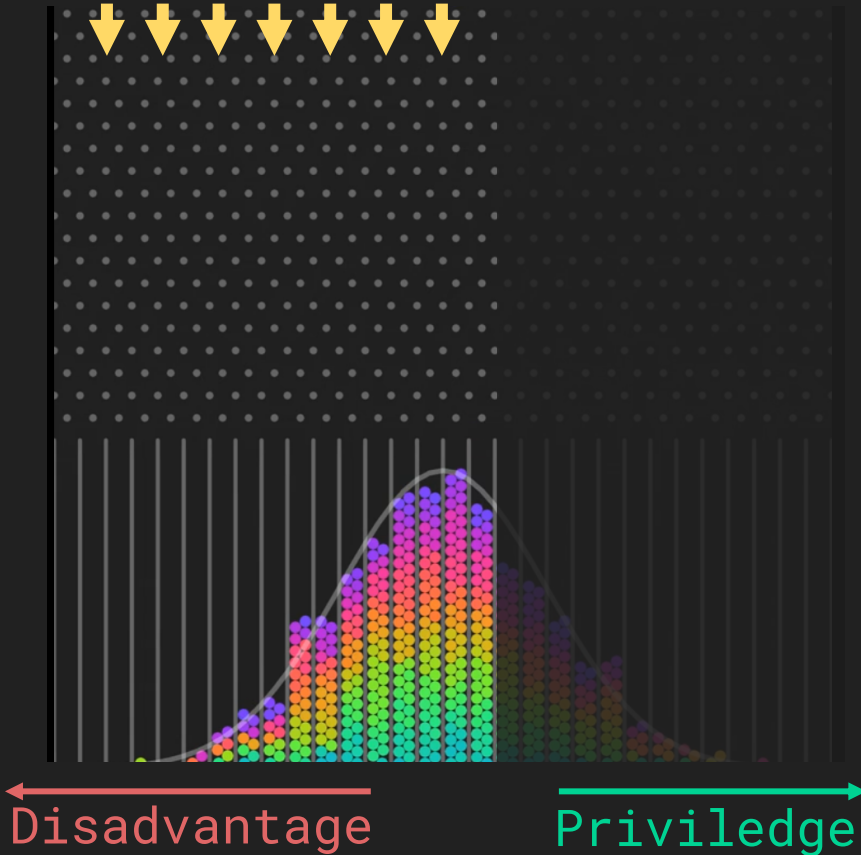


Education is one of the main ways we have to help disadvantaged populations.

Most educational opportunities are gated, which tends to help those that are ahead get farther ahead.



Why SEARCH?



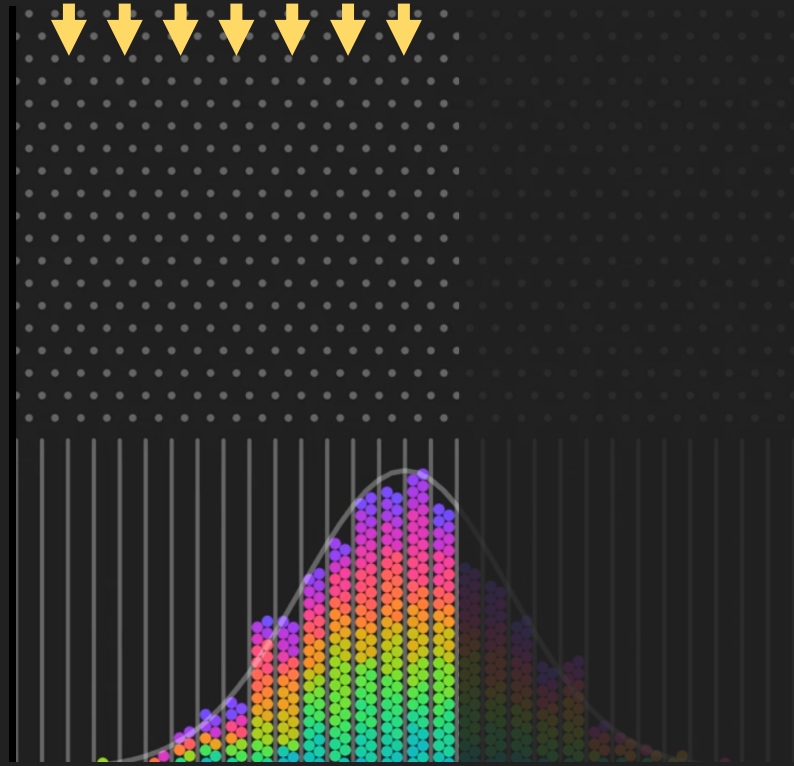
Education is one of the main ways we have to help disadvantaged populations.

Most educational opportunities are gated, which tends to help those that are ahead get farther ahead.

So we restrict them, and it helps, up to a point.



Why SEARCH?



Education is one of the main ways we have to help disadvantaged populations.

Most educational opportunities are gated, which tends to help those that are ahead get farther ahead.

So we restrict them, and it helps, up to a point.

Because privilege is relative.

advantage Priviledge



Sources of privilege we wanted to offset



- Socio-economic status, wealth, and place of origin:
 - Limited access to educational opportunities.
 - Significant difference in skillsets.
 - Limited access to computational resources.
- Gender, race, ethnicity, religion, and sexual orientation:
 - Subconscious and conscious bias.
 - Harder access to opportunities.
- Age, employment, family responsibility, veteran status:
 - Scientific education is focused on students.
 - Limited time availability to learn.
 - Reduced access to new techniques and technologies.



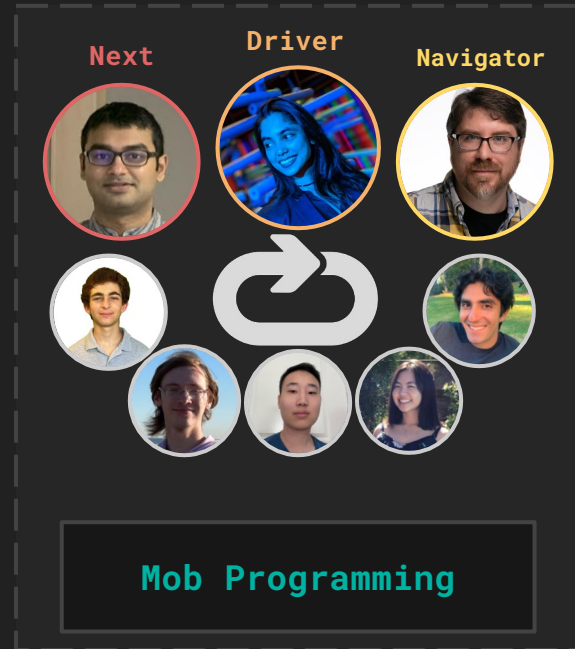
Program design



Significant difference in skillsets.
Reduced access to new techniques and technologies.

Collaborative programming (mob-coding)

- Equal sharing of development responsibility.
- Short, low-stakes rotations.
- Simultaneous development and documentation of code.
- Mob efficiently finds solutions to obstacles.
- Effective learning environment for equalizing skillsets.

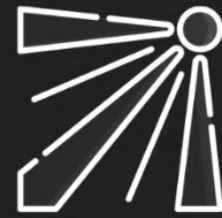




Program design



<p>Significant difference in skillsets. Reduced access to new techniques and technologies.</p>	Collaborative programming (mob-coding)
<p>Limited time availability to learn</p>	4 hours/week commitment, no assignments
<p>Limited access to computational resources.</p>	Development uses Google colab and visual studio code liveshare sessions
<p>Limited access to educational opportunities. Harder access to opportunities. Scientific education is focused on students.</p>	Anyone can join as long as they show up and actively engage in collaborative programming.



Subhamoy C.



Jasper D.



Maxwell R.



Cameron W.



Matin Q.



Daniel G.



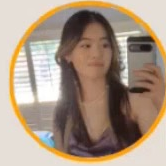
Nadia A.



Sierra M.



Jonathan V.



Jennifer L.



Scott M.



Ilya K.



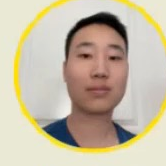
Spencer G.



Andres M.



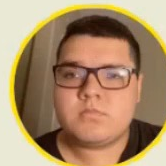
Julie C.



Justin G.



Jacob C.



Miguel T.



David S.



Opening: SEARCH objectives

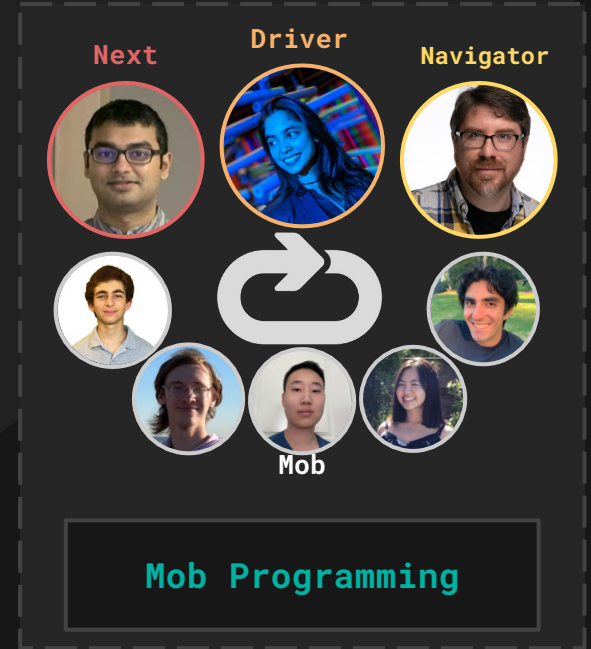
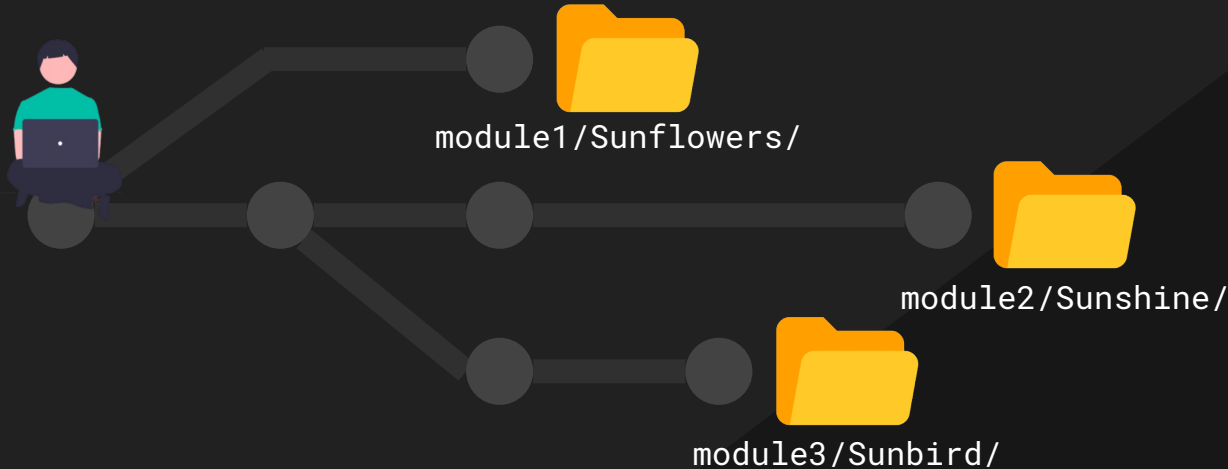


- ❑ To perform cutting edge **Machine Learning** (ML) research in **heliophysics** involving a diverse group of people without requiring *any prior ML technical expertise*.
- ❑ To create a space where **all are welcome** and where all participants **level up** and **work** together towards common technical goals.



Opening: Team Organization and Workflow

- VSCode Liveshare
- Mob programming
- Google Colab
- Team member rotation
- Github
- Weekly Code review
- Agile workflow

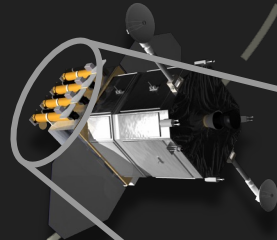
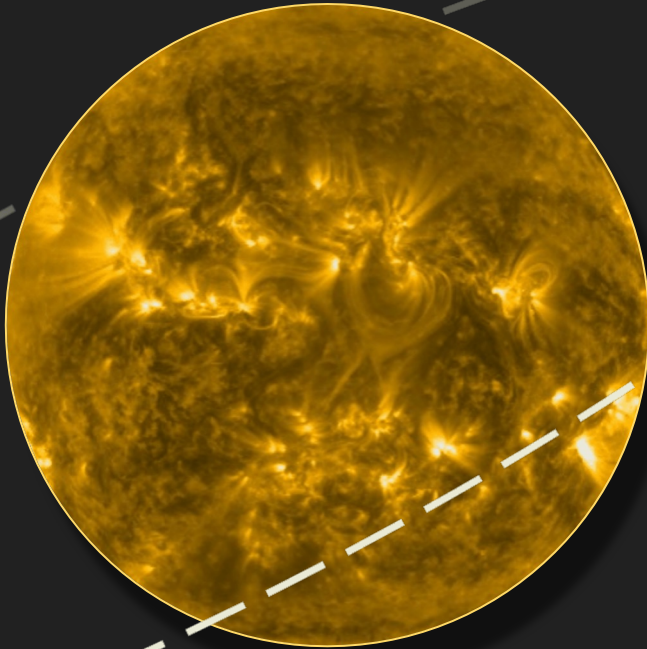




Introduction: Solar Dynamics Observatory



AIA (Atmospheric Imaging Assembly)

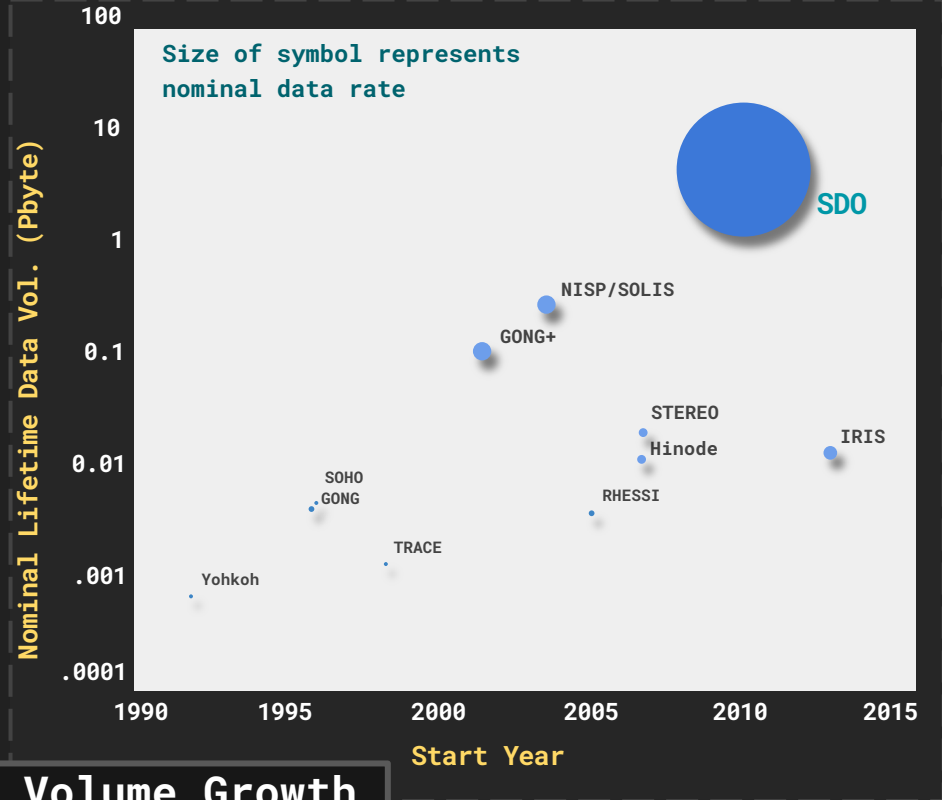




Introduction: AIA Data Volume & Motivation



- **Petabytes** of data
- Data must be **manually searched** through



Solar Data Volume Growth



Program Phases



EDUCATION

DATA

ML

USER INTERFACE





Components: Overall Pipeline



SDO Downloader

```
sdate: (str)
edate: (str)
instrument: (str)
wavelength: (list)
cadence: (str)
grayscale: (bool)
multiwavelength: (bool)
```

Stanford Joint
Science
Operations
Center (JSOC)



Solar Images



Dataset (AIA)

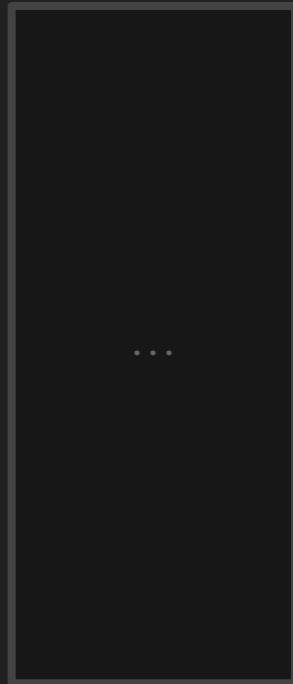
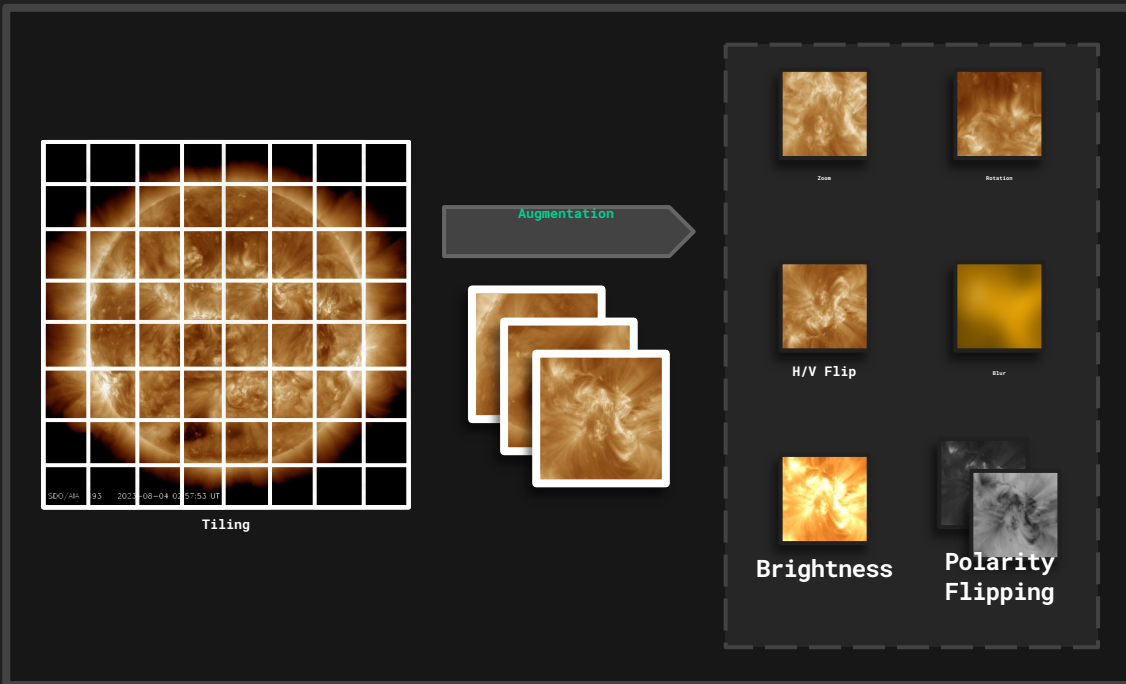
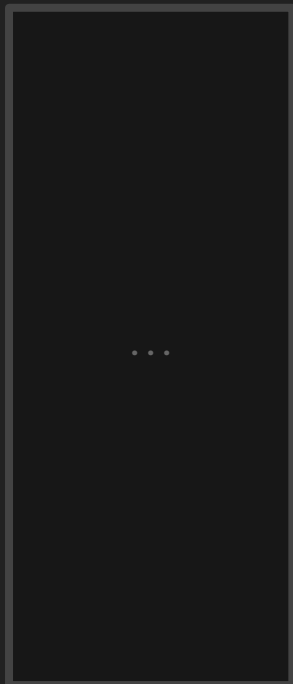
Solar Dynamics Observatory Downloader

AI/ML Ready
Packager

Similarity
Search Engine



Components: Overall Pipeline



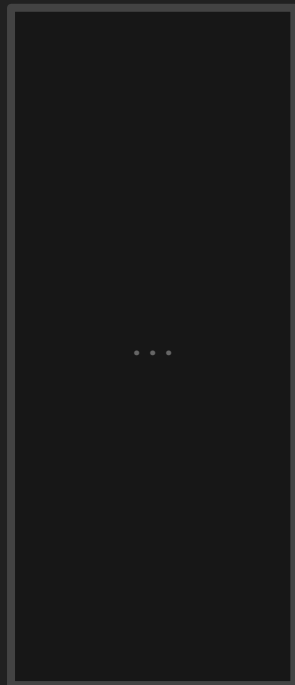
SDO
Downloader

AI/ML Ready Packager

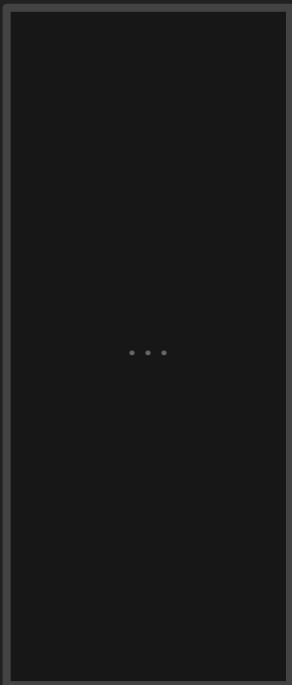
Similarity
Search Engine



Components: Overall Pipeline



SDO
Downloader



AI/ML Ready
Packager



Similarity Search Engine



Introduction to Self-Supervised Learning: General Overview (Model-Agnostic)



Trained by comparing:

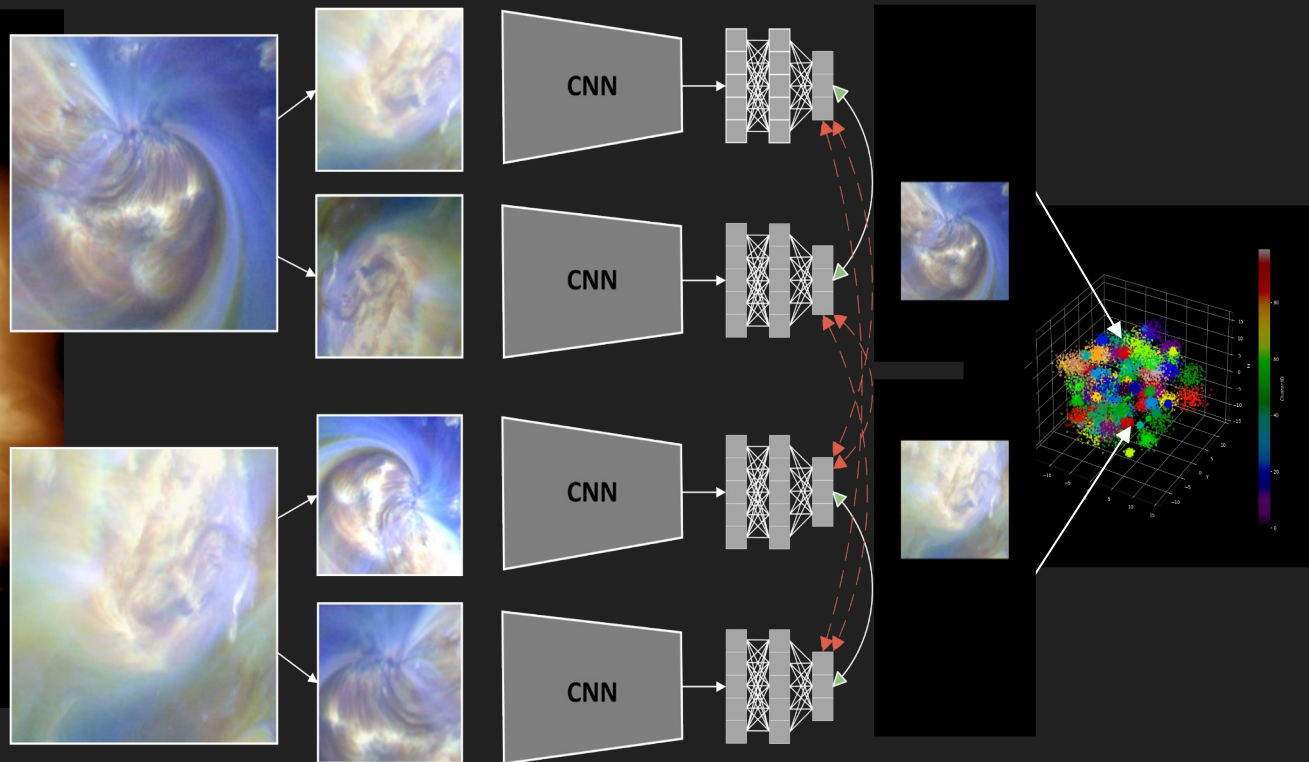
Image vs. transformed image (and/or a contrasting image)

Allows us to discern similar & dissimilar images.

Returns:

Similar embeddings for an image & its transformed counterpart.

For dissimilar images the embeddings should be very different.





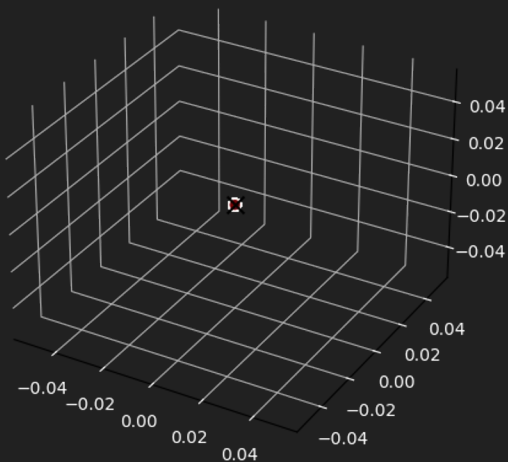
Introduction to Self-Supervised Learning: General Overview (Model-Agnostic)



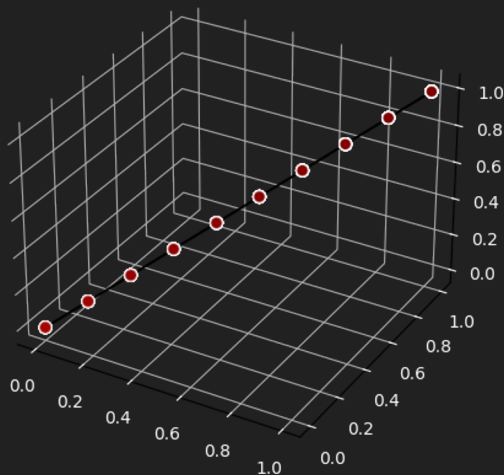
Collapse:

Solution collapse can occur when a model fails to learn meaningful feature representations from the data and converges on a single embedding or limited embedding space.

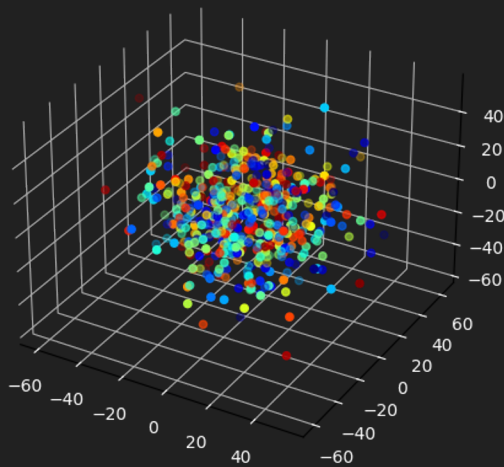
Complete Collapse



Dimensional Collapse



No Collapse

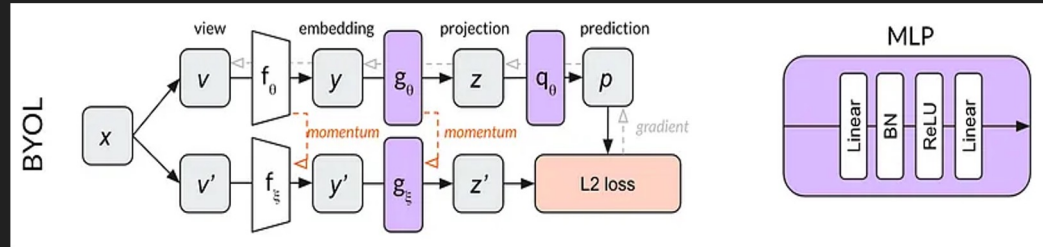
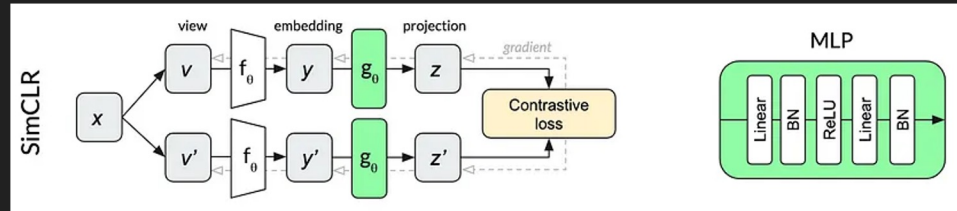
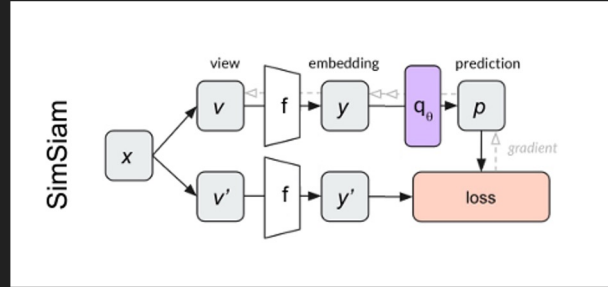




Introduction to Self-Supervised Learning: General Overview (Model-Agnostic)

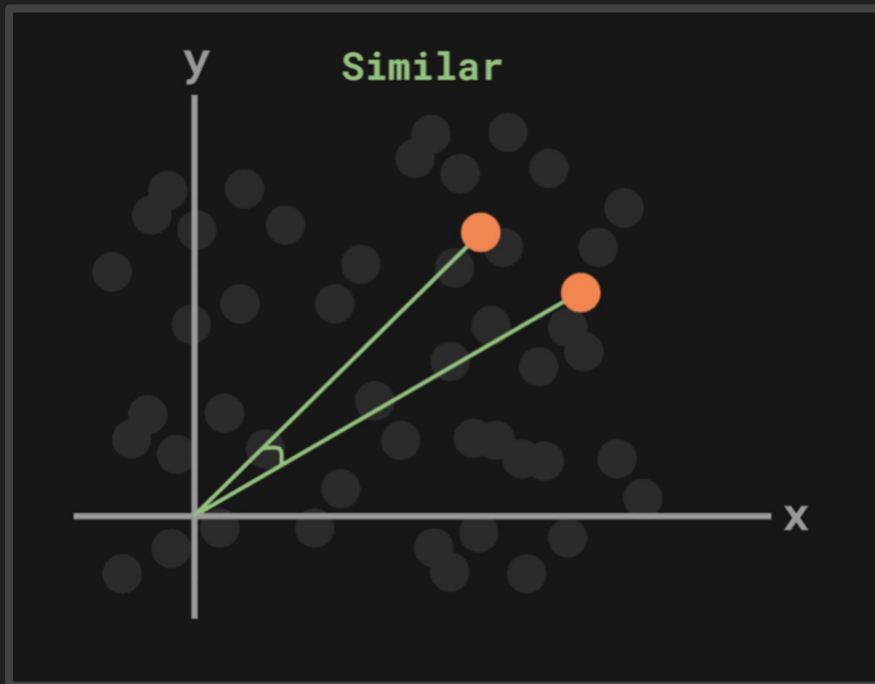


ML algorithms we used:
SimSiam; SimCLR; BYOL
("Bootstrap Your Own Latent").

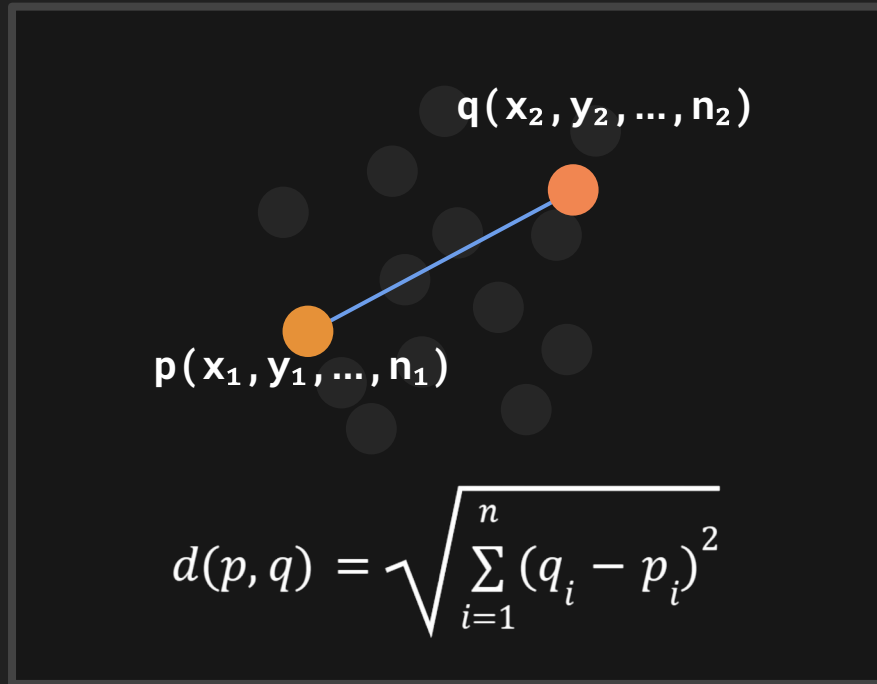




Introduction to Self Supervised Learning: Metrics and Losses



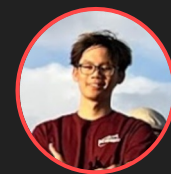
Cosine Similarity



Euclidean Distance



Introduction to Self Supervised Learning: Metrics and Losses

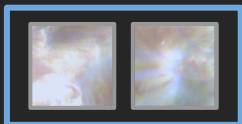


Contrastive Loss

Measures how well the model can **contrast** between **similar** and **dissimilar** data points

1) Augmented Image in Batch

Pair 1

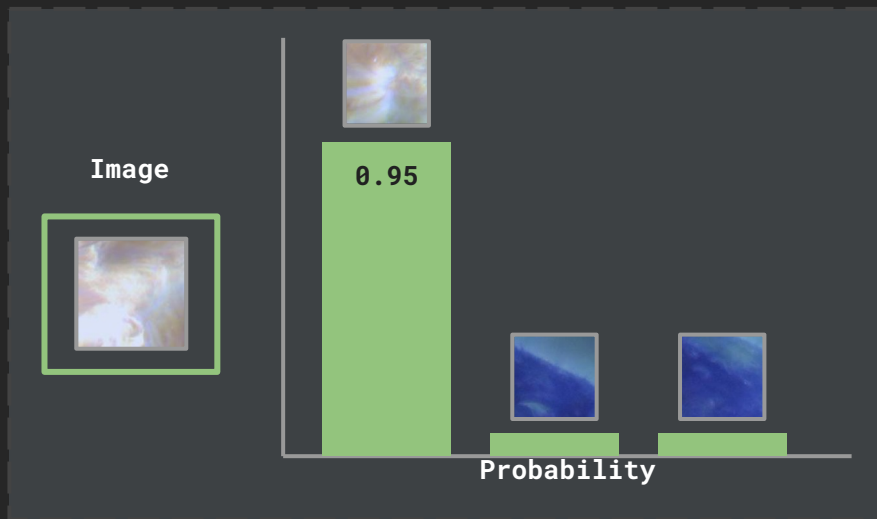


Pair 2



2)

$$\text{Softmax} = \frac{e^{\text{similarity}(\text{img}_1, \text{img}_1)}}{e^{\text{similarity}(\text{img}_1, \text{img}_1)} + e^{\text{similarity}(\text{img}_1, \text{img}_2)} + e^{\text{similarity}(\text{img}_1, \text{img}_3)}}$$





Introduction to Self Supervised Learning: Metrics and Losses



3) Noise Contrastive Estimation (NCE)

$$l(\text{img}_1, \text{img}_2) = -\log(\text{Softmax}(\text{img}_1, \text{img}_2))$$

$$l(\text{img}_3, \text{img}_4) = -\log(\text{Softmax}(\text{img}_3, \text{img}_4))$$

$$l(i, j) = -\log\left(\frac{\exp(s_{ij})}{\sum_{k=1}^N \exp(s_{ik})}\right)$$

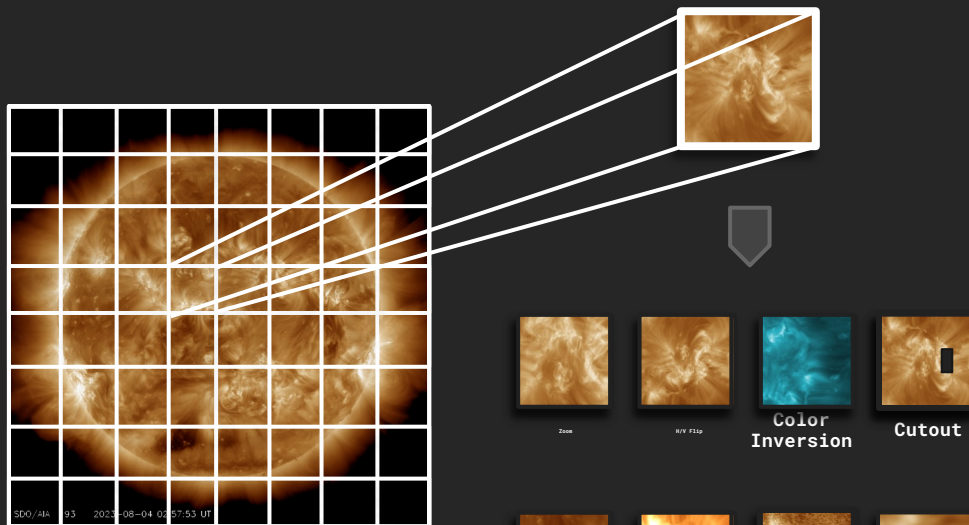
4) Compute Loss

$$L = \frac{[l(\text{img}_1, \text{img}_2) + l(\text{img}_3, \text{img}_4)] + [l(\text{img}_5, \text{img}_6) + l(\text{img}_7, \text{img}_8)]}{2 * N}$$

$$L = \frac{1}{2N} \sum_{k=1}^N [l(2k - 1, 2k) + l(2k, 2k - 1)]$$



Introduction to Self Supervised Learning: Augmentations



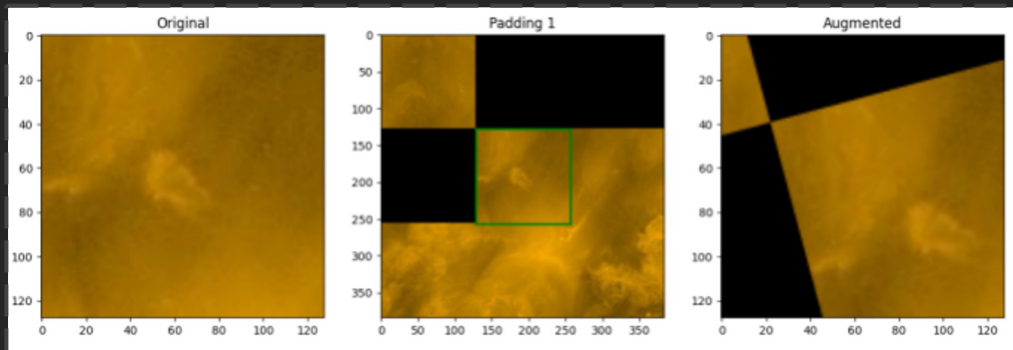
Super Image

Augmenting Data

- ✓ Introduce **Variations**
- ✓ Robustness to **Noise**
- ✓ **Generalizability**
- ✓ Prevents Data **Scarcity**
- ✓ Reduces **Overfitting**



Introduction to Self Supervised Learning: Augmentations

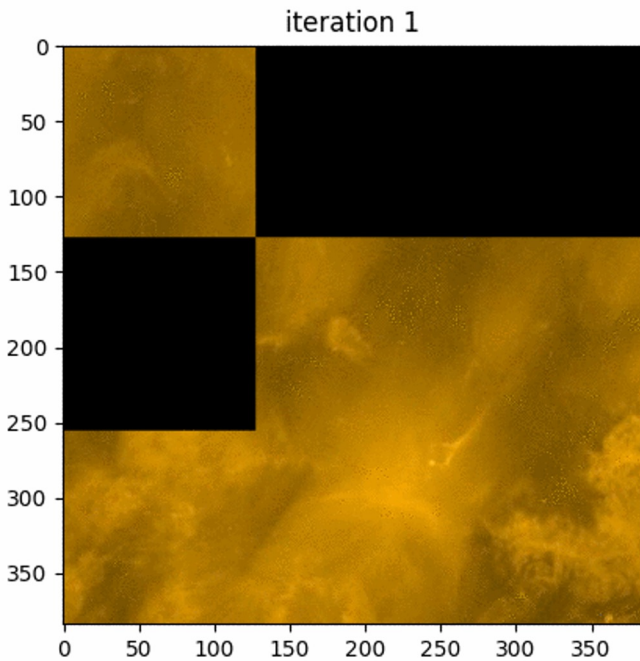


Augmenting Data

- ✓ Augmentations, and train/val split creates gaps
- ✓ Nearest neighbor VS Iterative filling
- ✓ Iterative filling creates improvement in terms of collapse



Introduction to Self Supervised Learning: Augmentations

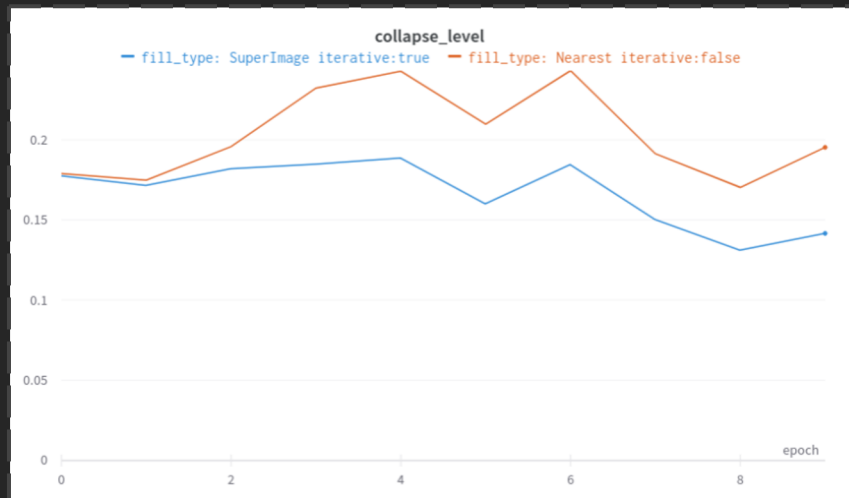
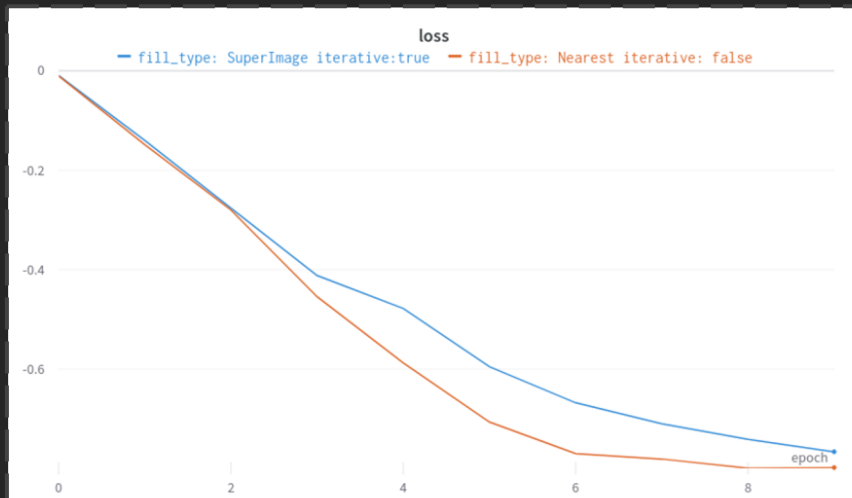


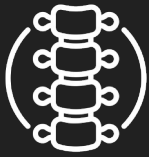
Augmenting Data

- ✓ Augmentations, and train/val split creates gaps
- ✓ Nearest neighbor VS Iterative filling
- ✓ Iterative filling creates improvement in terms of collapse

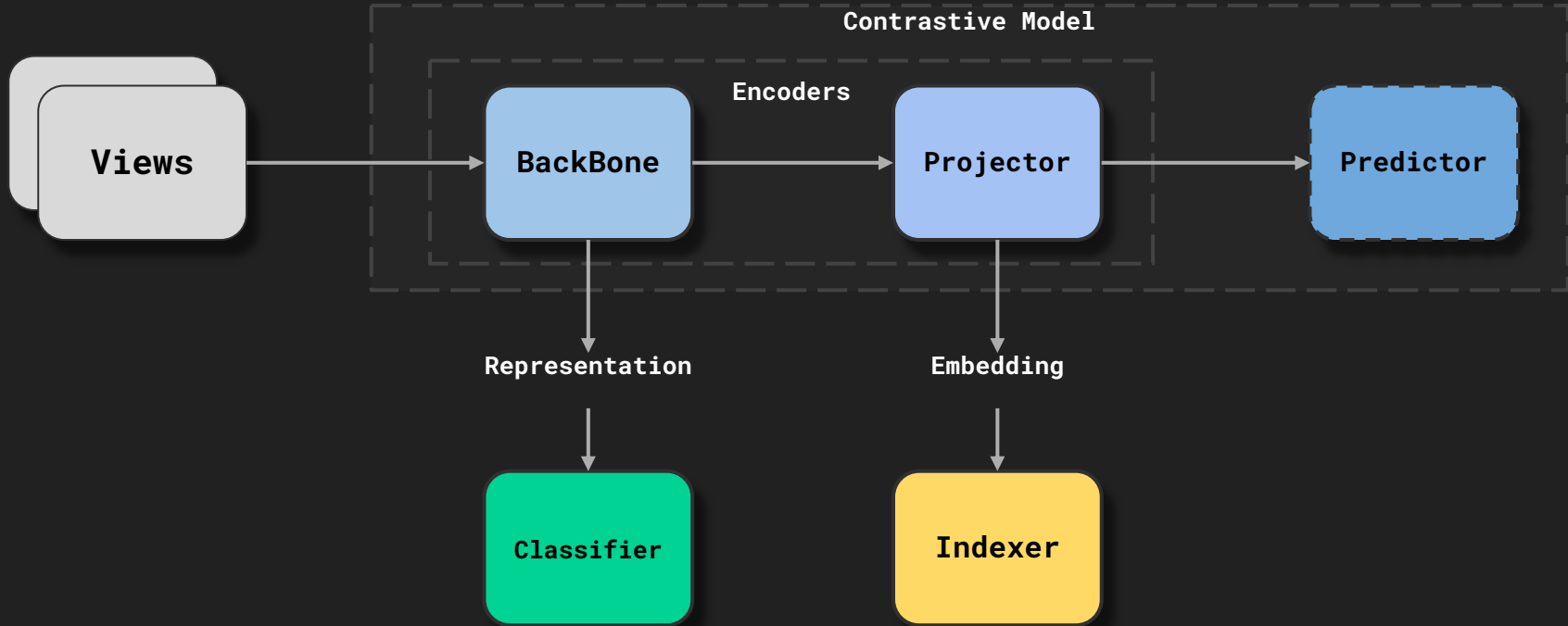


Introduction to Self Supervised Learning: Augmentations





Introduction to Self Supervised Learning: Embeddings from backbone vs. projection head





Introduction to Self Supervised Learning: Embeddings from backbone vs. projection head



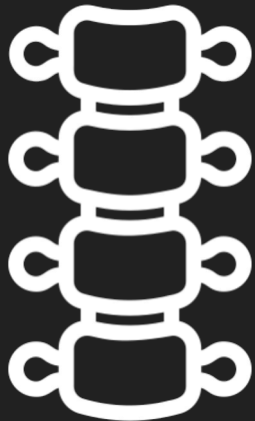
SEARCH
Projection
Head



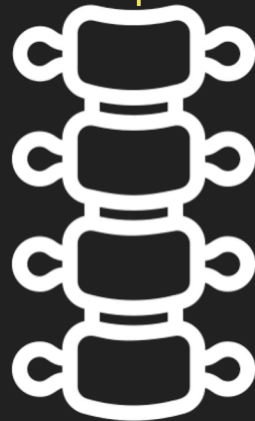
NEW
Projection
Head



SEARCH
Backbone

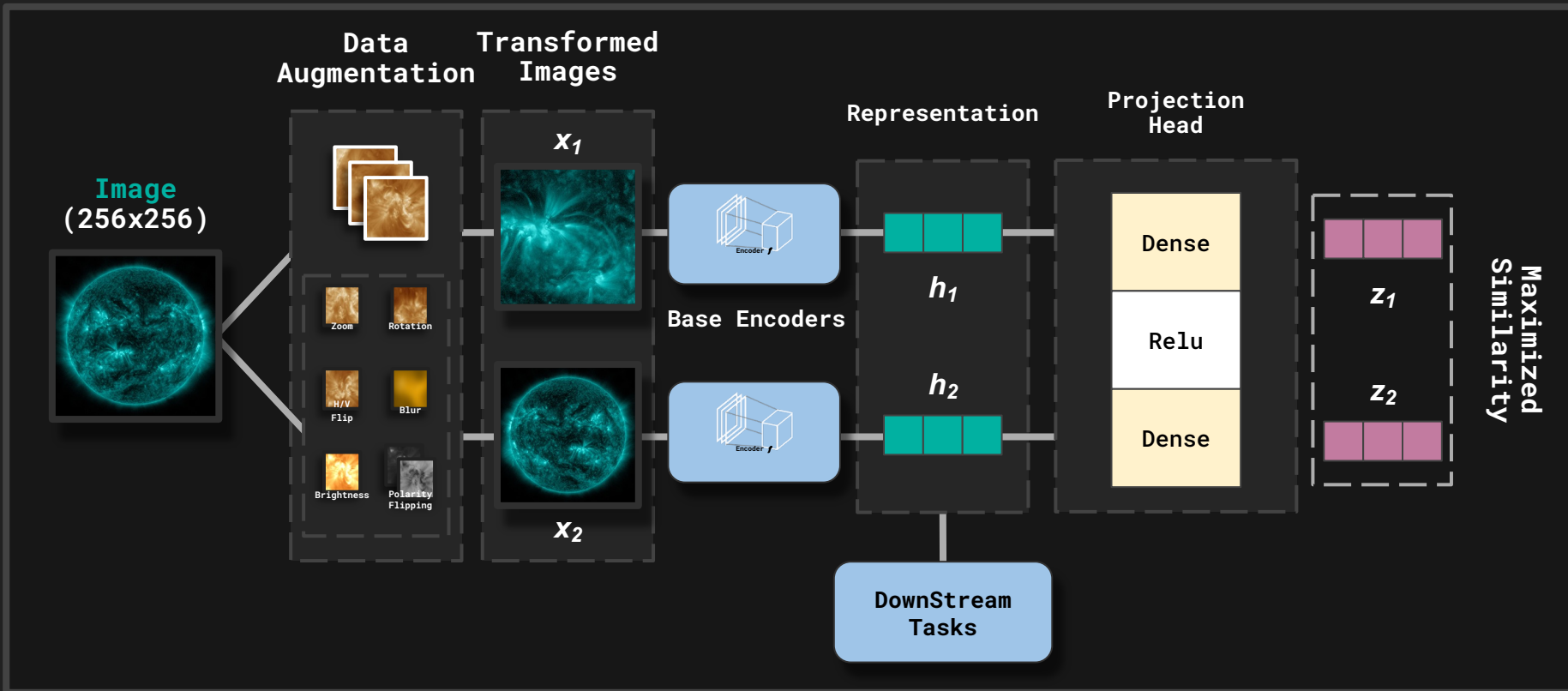
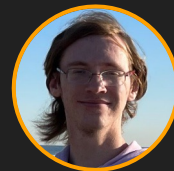


SEARCH
Backbone



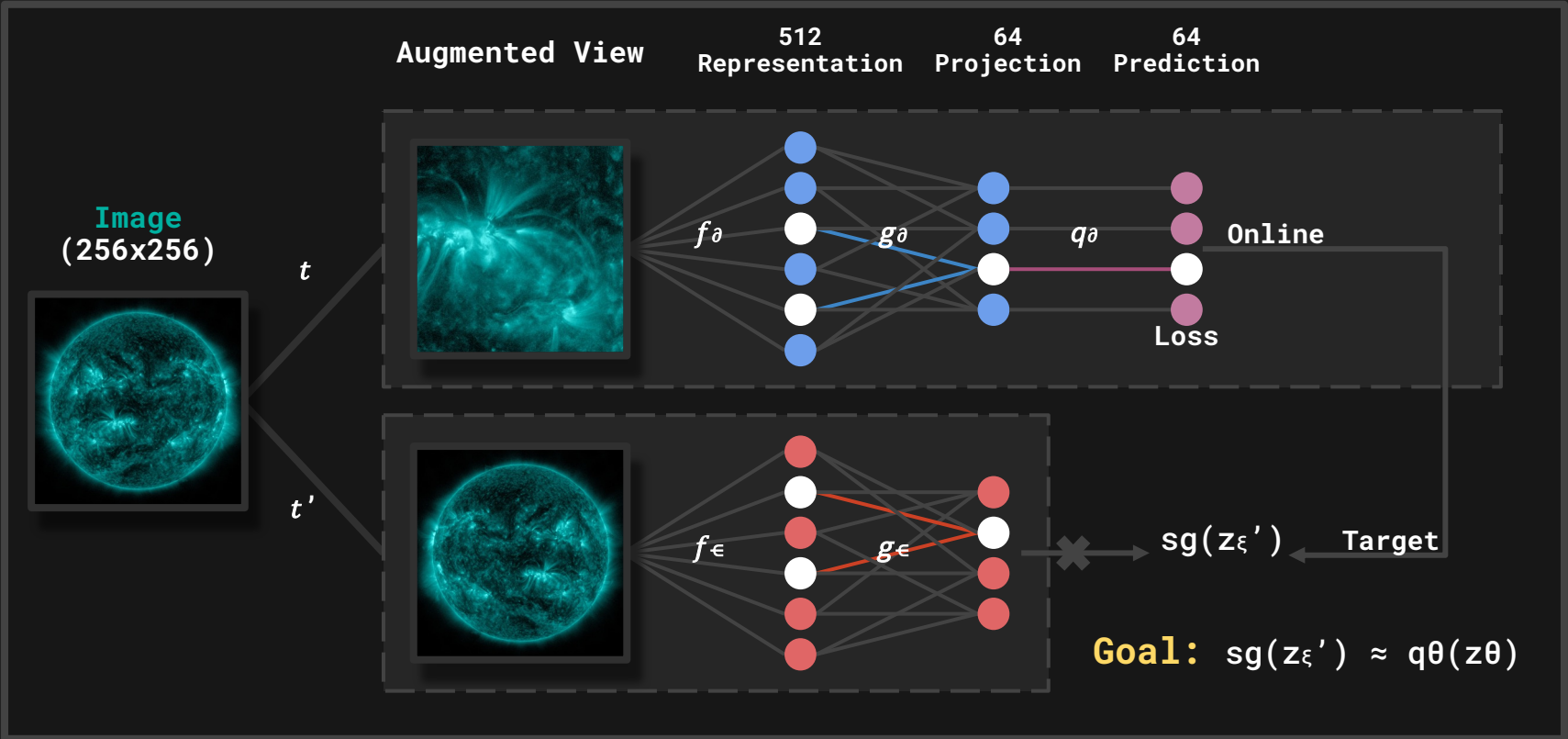


Introduction to Self Supervised Learning: SimCLR



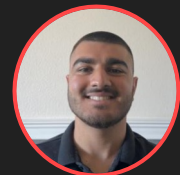


Introduction to Self Supervised Learning: BYOL: Bootstrap Your Own Latent





Introduction to Self Supervised Learning: SimSiam



SimSiam

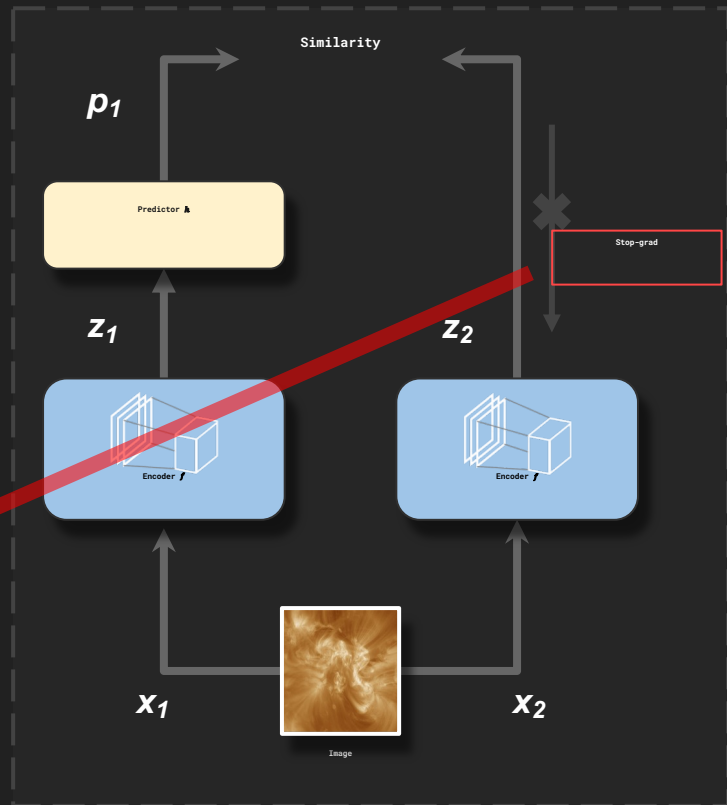
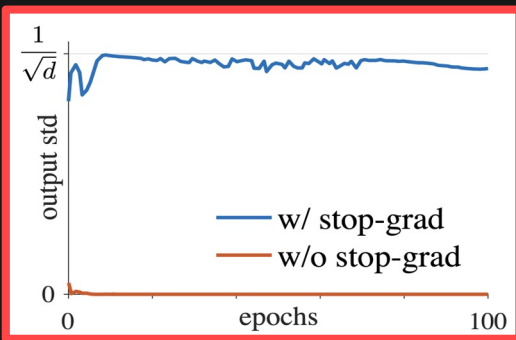
- ✓ Identical encoders
- ✓ Stop-gradient operation is sufficient to prevent collapse
- ✓ Back-propagation through only one branch

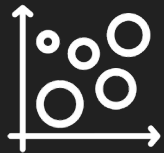
+ Stop-gradient operation is critical for preventing collapse

+ Collapse level measured through standard deviation of vectors.

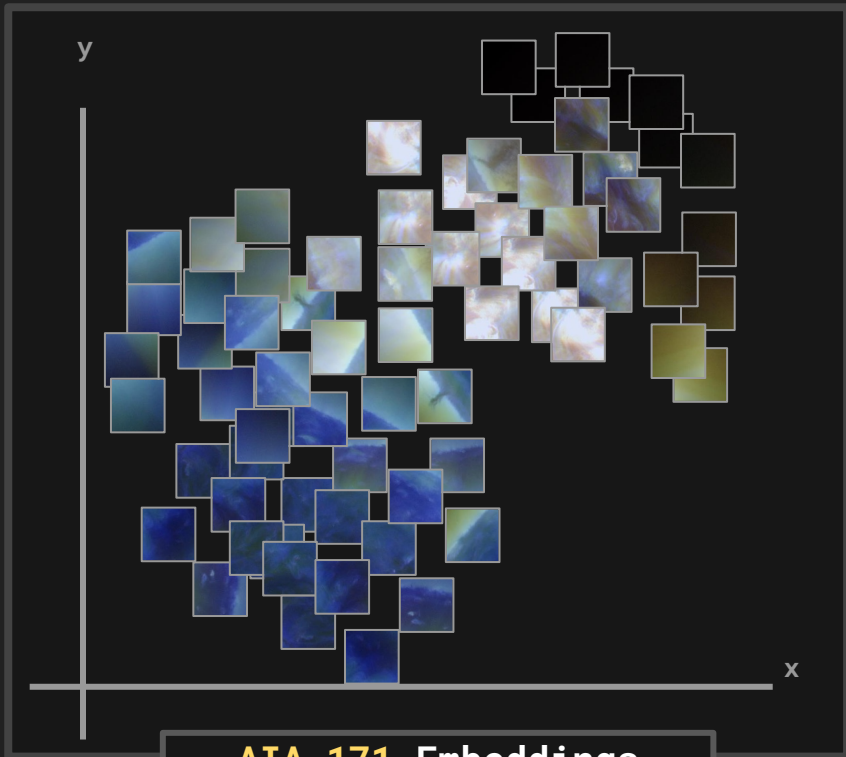
BYOL

- ✗ Target branch is created as a moving average of online branch
- ✗ Utilize various architectural asymmetries to prevent collapse
- ✓ Back-propagation through only one branch

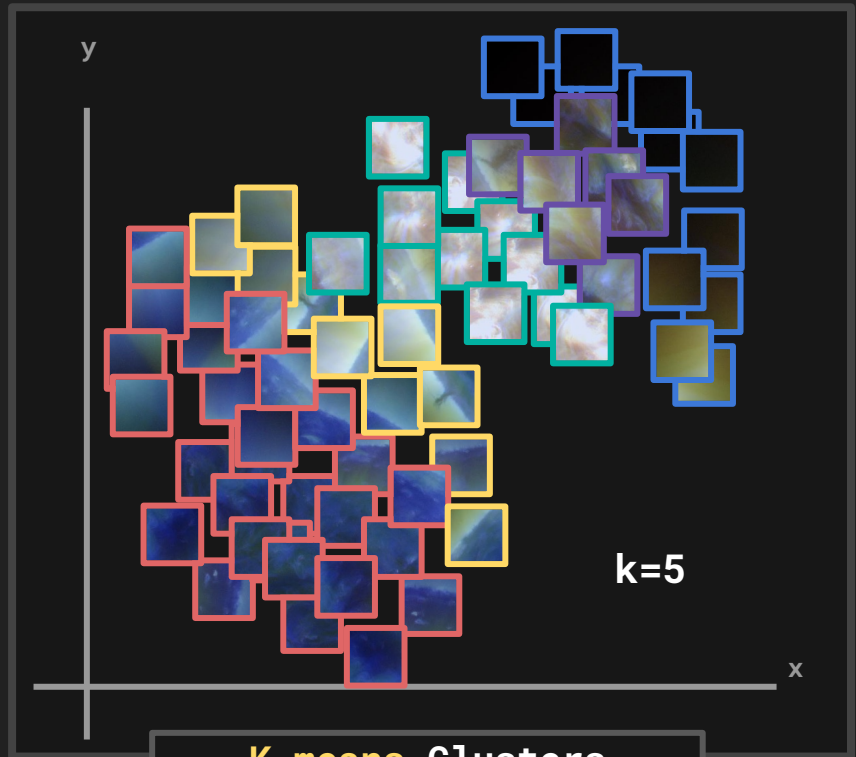




Introduction to Self Supervised Learning: Clustering



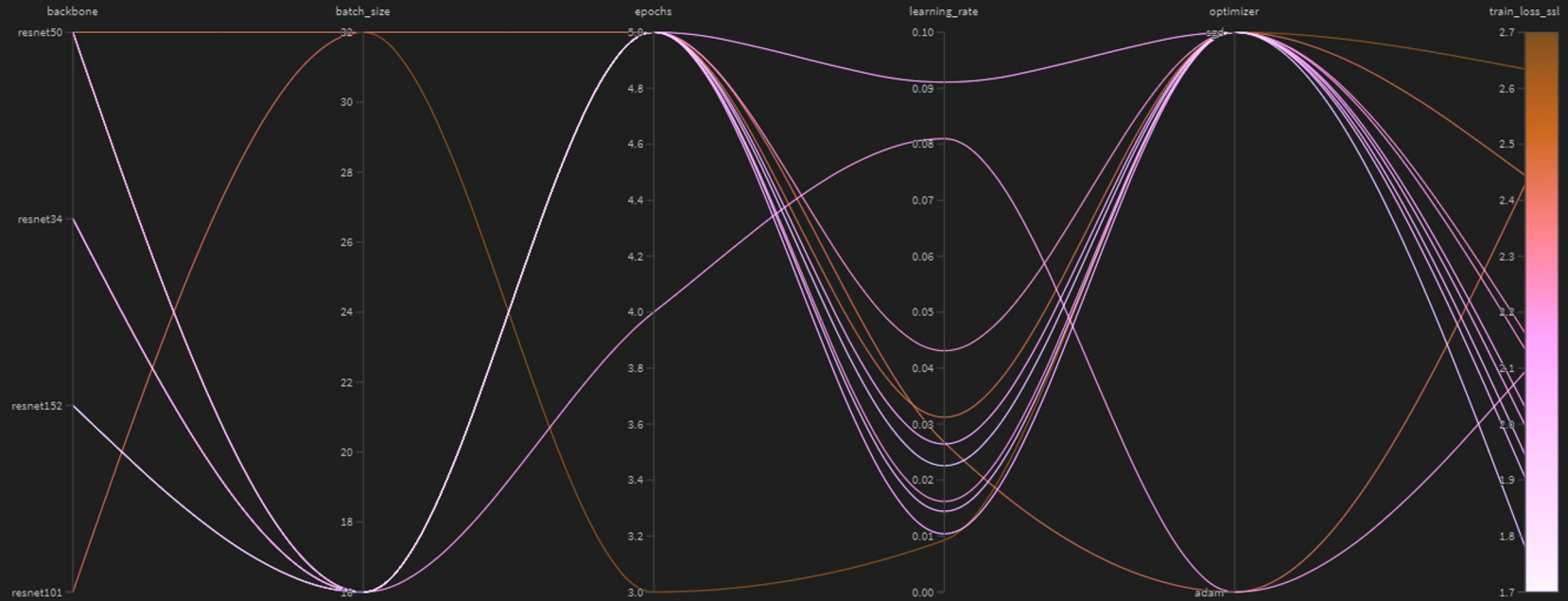
AIA 171 Embeddings



K-means Clusters



SimCLR: Sweep

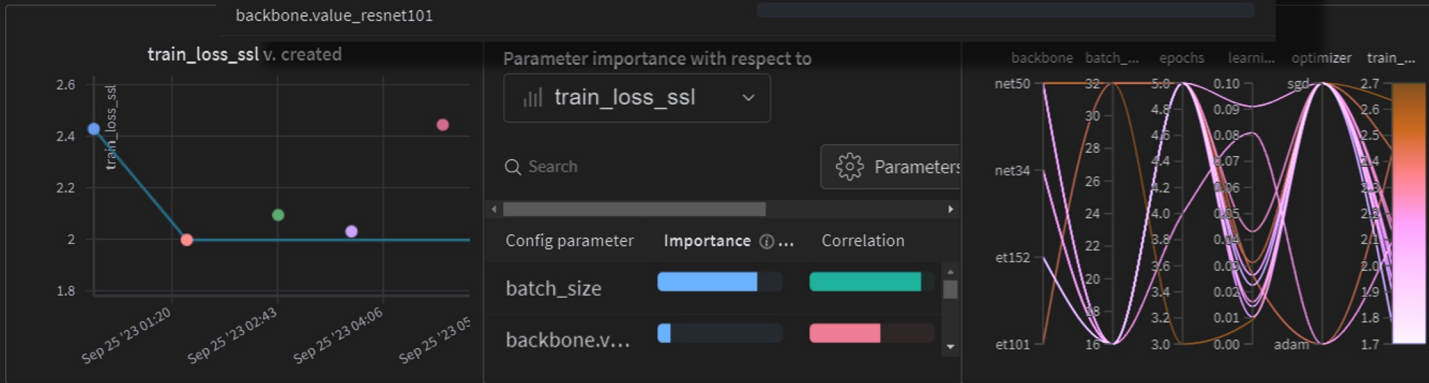
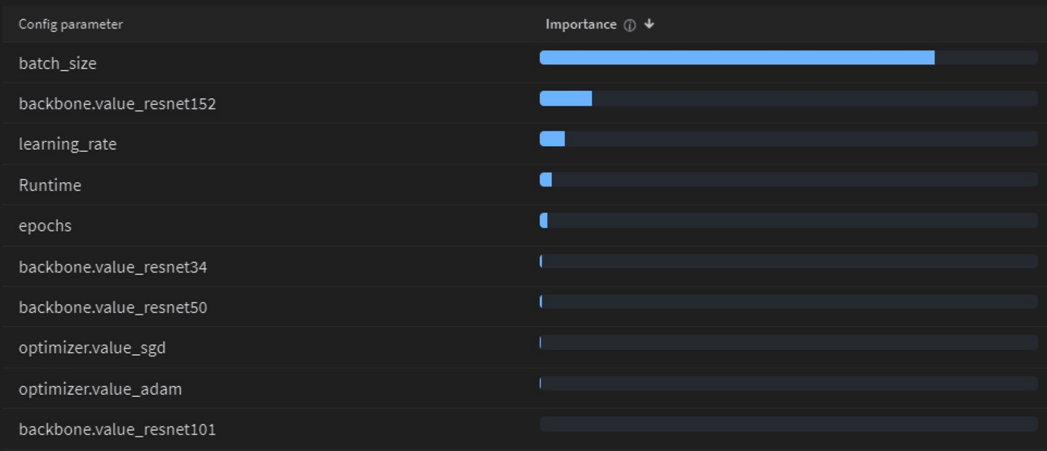
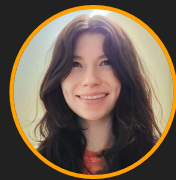


Features:

- "NT-Xent loss" (Normalized Temperature-Scaled Cross-Entropy Loss)
- Additional augmentations

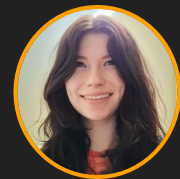


SimCLR: Hyperparameters

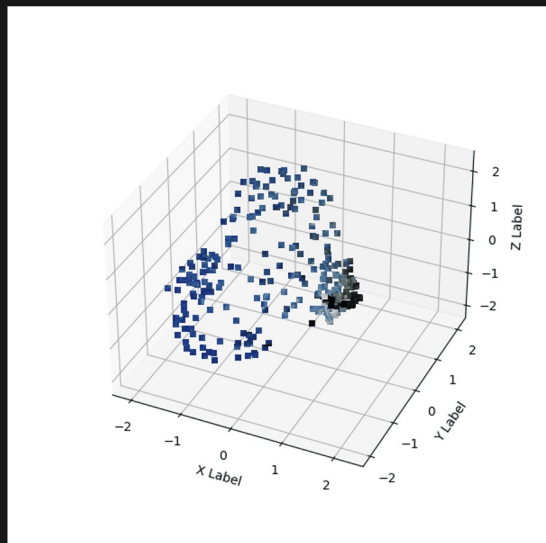




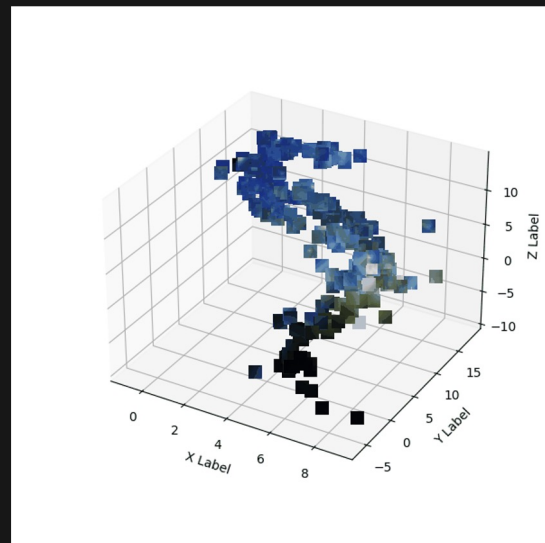
SimCLR: Plots



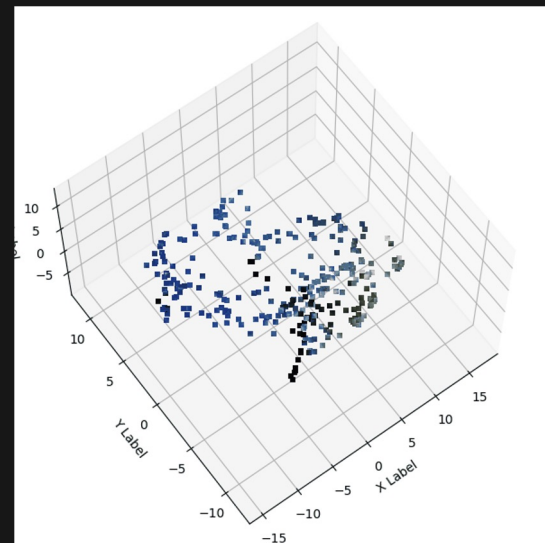
PCA



UMAP

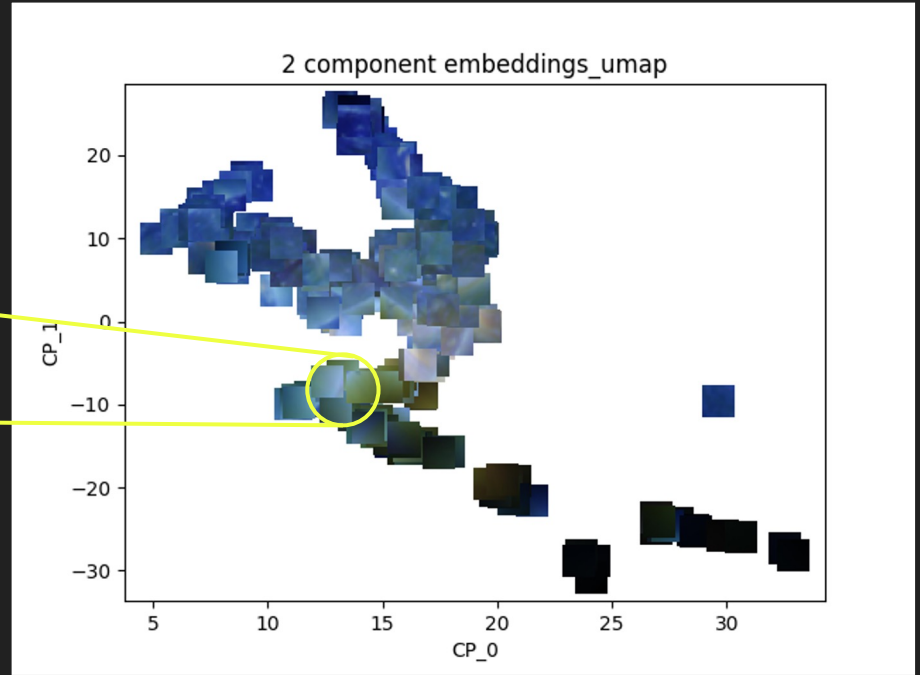
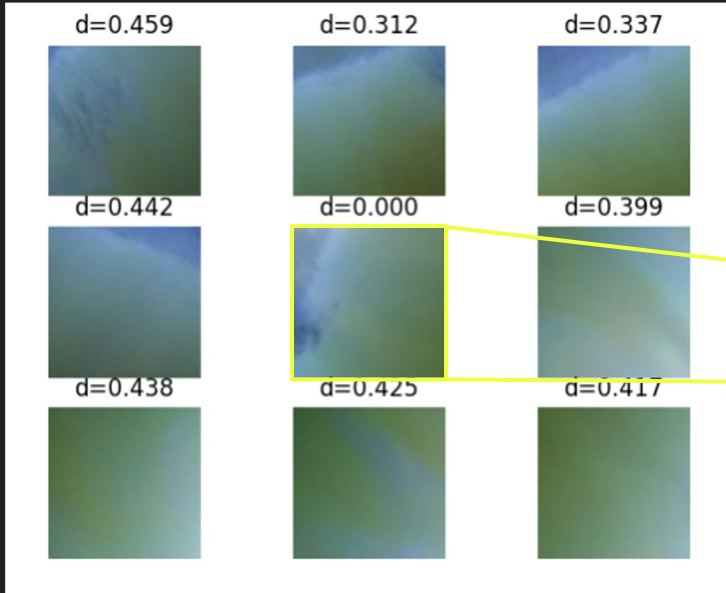
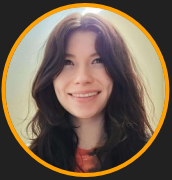


TSNE



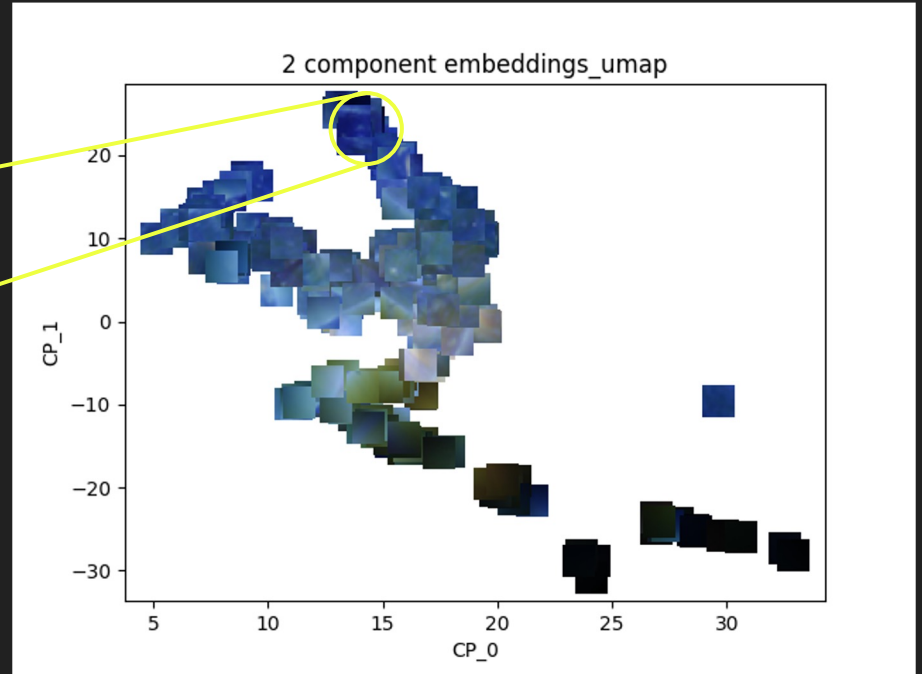
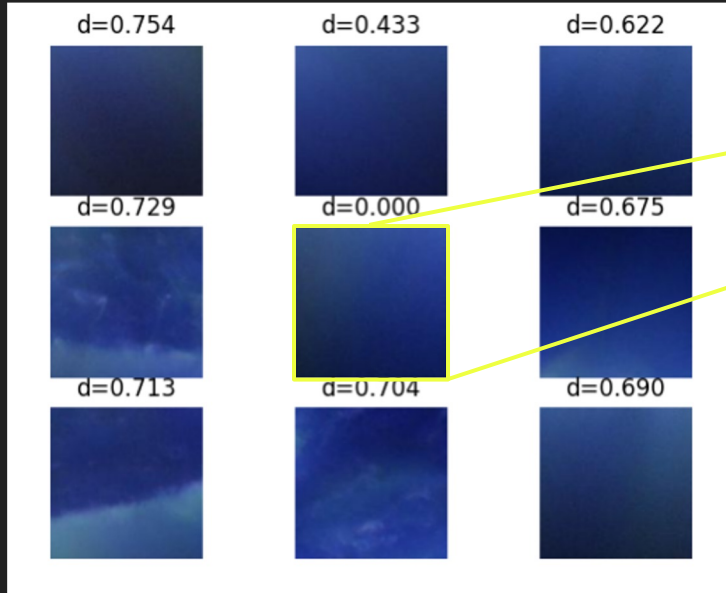
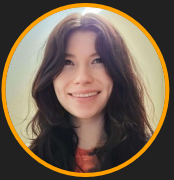


SimCLR: K-Nearest Neighbors



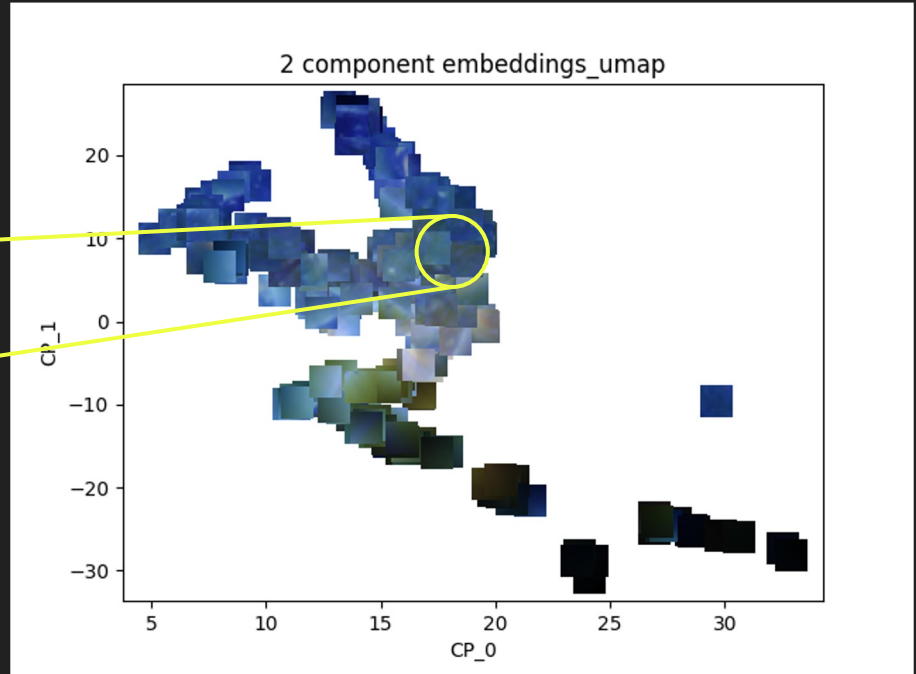
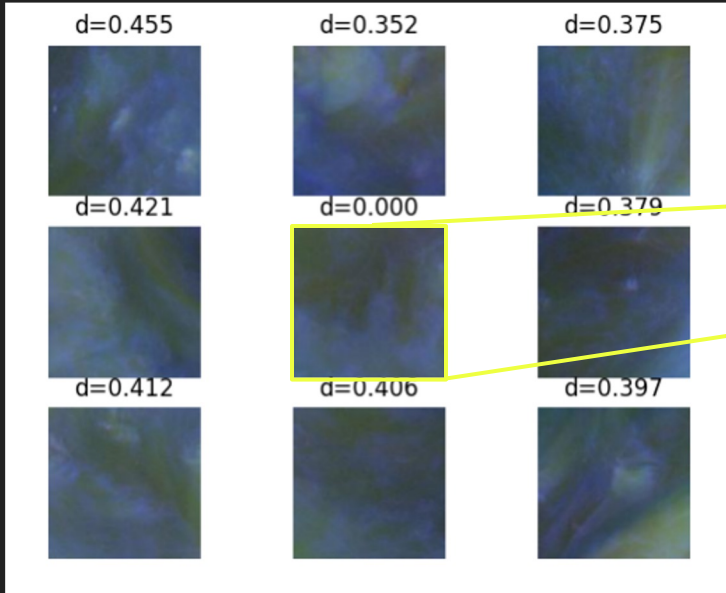
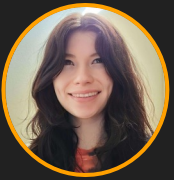


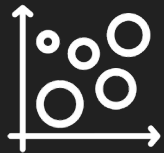
SimCLR: K-Nearest Neighbors



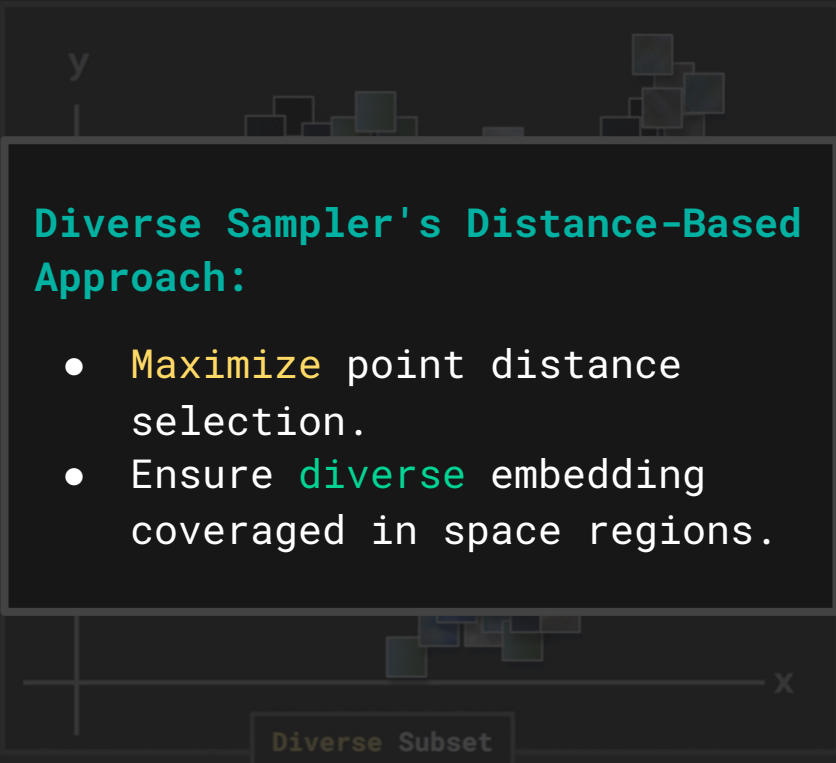
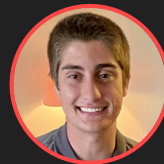


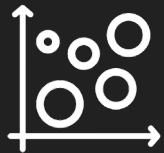
SimCLR: K-Nearest Neighbors



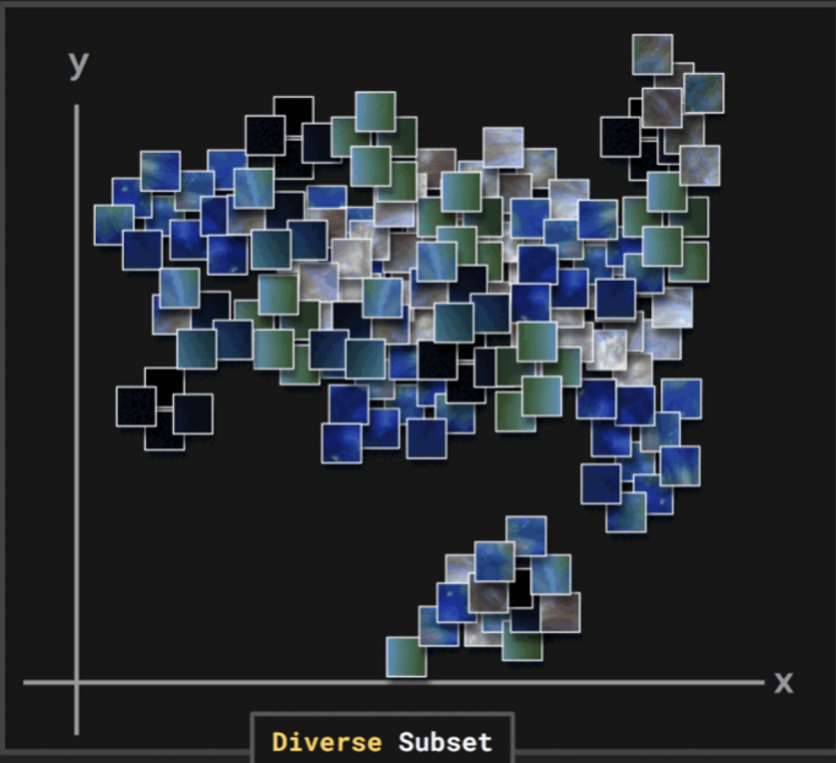
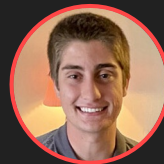


Team Highlights: Diverse sampling and Iterative training





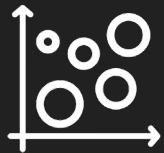
Team Highlights: Diverse sampling and Iterative training



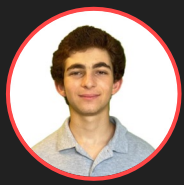
Enhanced Training with Diverse Sampling:

- Custom **training** loop integration.
- Iteratively **update** dataset with **diverse** samples.
- **Broaden** model exposure to capture data variations.

Full Subset



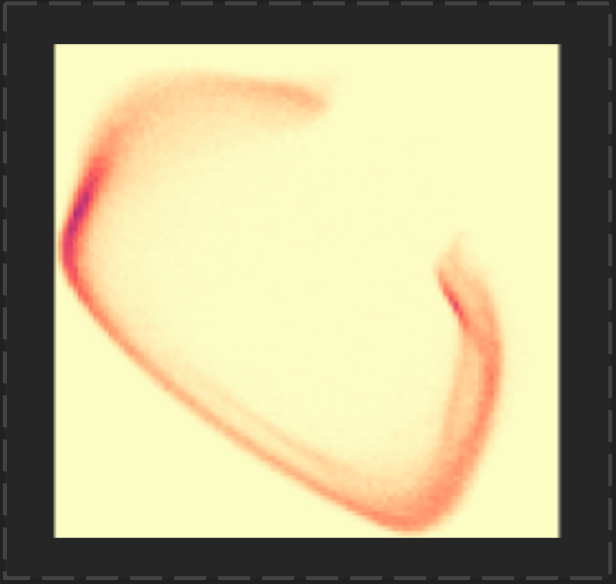
Team Highlights: Contrastive vs. Non-contrastive losses on embedding space



2D Histogram of Embedding Spaces



Non-Contrastive Loss



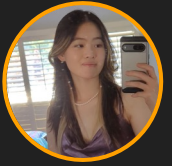
Contrastive Loss



Self-Supervised Search Demos
involving a User Query



Concluding Remarks: Implications and Potential



Collaborative Learning



Automation



Advances in Heliophysics



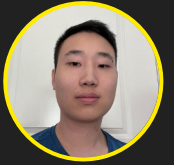
Concluding Remarks: Future Work



- ❑ Fine-Tuning
 - ❑ Using *Labeled Data*
 - ❑ Different *Algorithms and Architecture*
- ❑ Improving JSOC **Usability & Accessibility**
- ❑ **Optimization & Expansion**
- ❑ ML Democratization outside **Educational** Institutions & Open **Research** Accessibility
- ❑ Continue **Experimentation**



Concluding Remarks: Conclusion



- We have brought together passionate individuals from diverse backgrounds in **machine learning** and **heliophysics**
- From humble beginnings, we have built **advanced** data and machine learning systems that have made the intricate world of **scientific research** a bit more **approachable**
- We are also **proud** of what we have done and we look forward to accomplishing more in the **future**

