



# Unveiling Bird Morphology Evolution through Deep-Learning Based Image Embedding and Gene Association Testing

Final presentation gene2birds group A

Joshua Günther

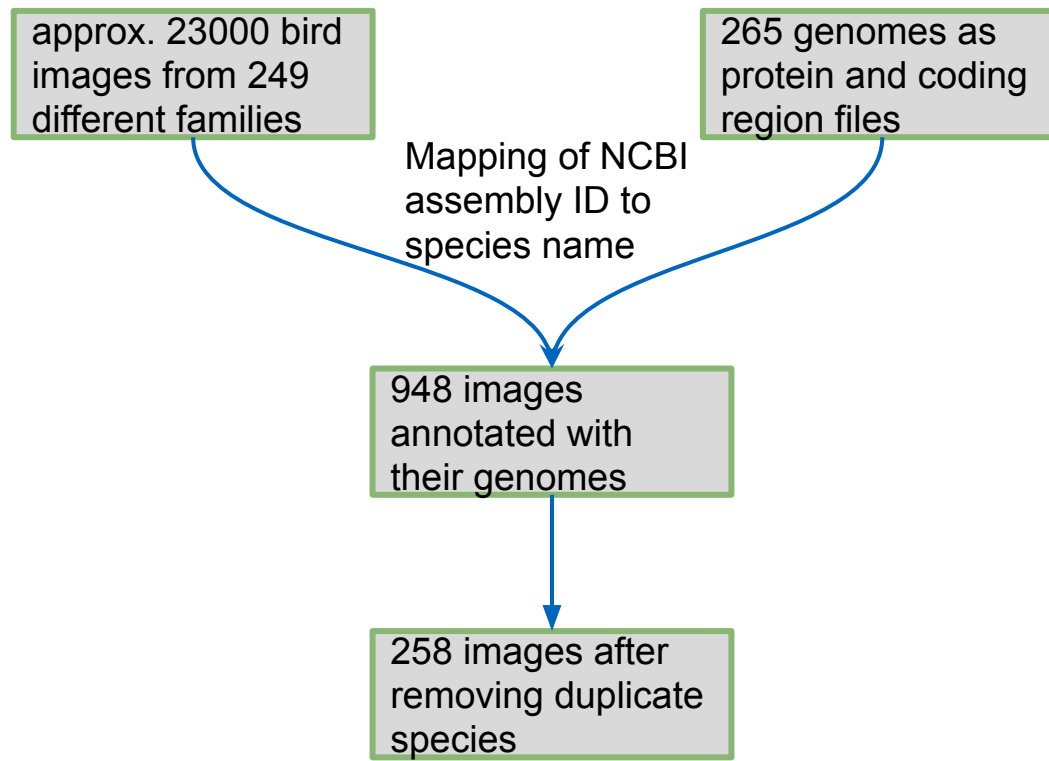
Berfin Erdoğan

Ufuk Demir

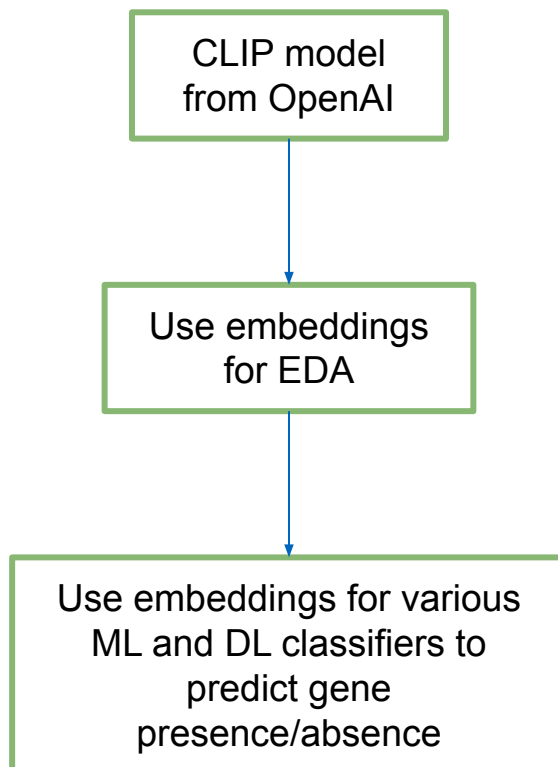
Mahima Arunkumar



# Input data



# Our approach



# Absence/presence table by species

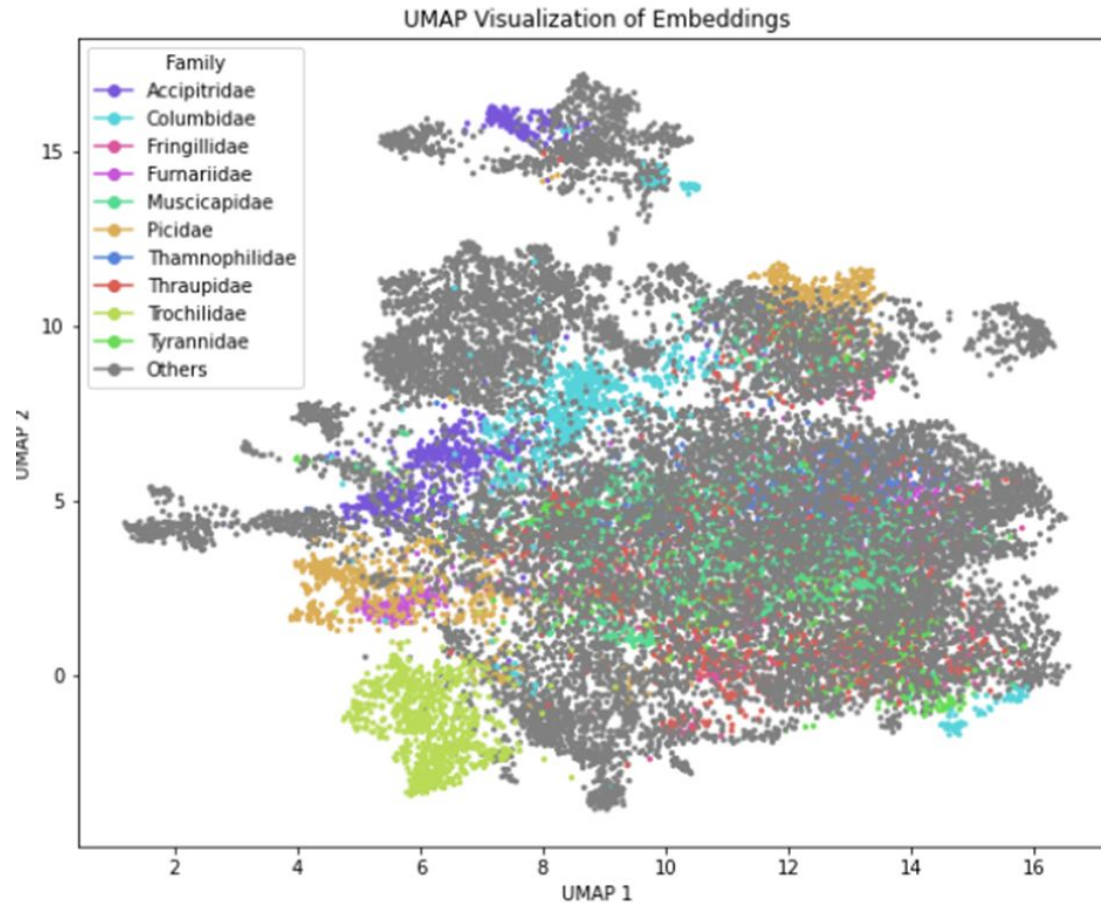
	Nothoprocta ornata	Smithornis capensis	Formicarius rufipectus	Sylvia atricapilla	Lanius ludovicianus	Amazona gouldingii	Probosciger atterrimus	Eolophus roseicapilla	Chunga burmeisteri	Herpetotheres cachinnans
<b>104K protein</b>	False	False	False	False	False	False	False	False	False	False
<b>110KD protein</b>	False	False	False	False	False	False	False	False	False	False
<b>1433B protein</b>	True	False	True	True	False	True	True	True	True	True
<b>1433E protein</b>	True	True	True	True	True	True	True	True	False	False
<b>1433F protein</b>	True	True	True	True	True	True	True	True	True	True
...	...	...	...	...	...	...	...	...	...	...
<b>ZY11B protein</b>	True	True	False	True	True	True	True	True	True	True
<b>ZYX protein</b>	False	True	True	False	True	True	True	True	True	False
<b>ZZEF1 protein</b>	True	True	False	True	True	True	True	False	True	True
<b>ZZZ3 protein</b>	True	True	False	True	False	True	False	False	True	True
<b>protein</b>	False	False	False	False	False	False	False	False	False	False

17060 rows × 258 columns

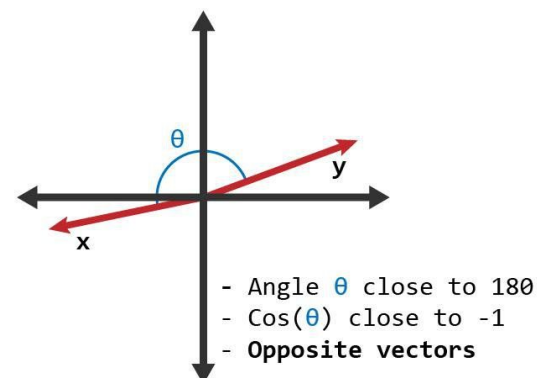
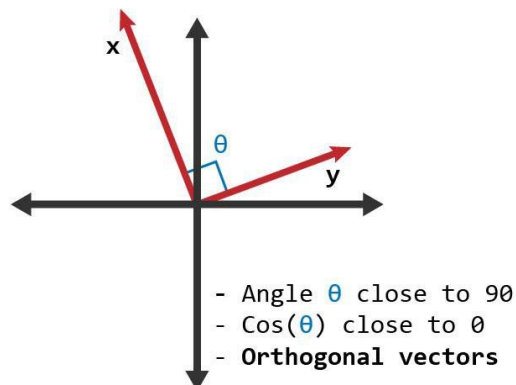
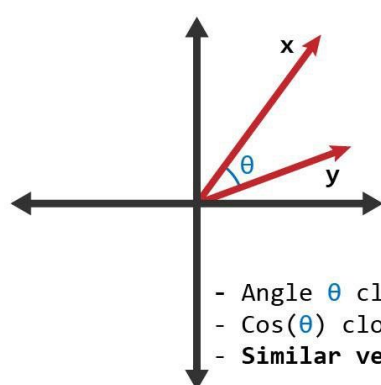
# Generated embeddings with CLIP

All 512 CLIP embedding features					
All 258 species	Nothoprocta ornata	0.4072875	0.36193112	-0.31036815	0.13646944
	Smithornis capensis	0.3131959	0.20072602	-0.43035	0.18086748
	Formicarius rufipectus	0.37125307	-0.050516717	-0.52256465	0.25555477
	Sylvia atricapilla	0.21032035	0.24915919	-0.6387599	0.5333118
	Lanius ludovicianus	0.6003466	0.18232825	0.022367803	0.04154638
	Amazona guildingii	0.33798376	0.14806776	-0.47427034	-0.19996244
	Probosciger aterrimus	0.36870858	-0.056442954	-0.33242458	-0.04306676
	Eolophus roseicapilla	0.40817946	0.23949404	-0.13612387	0.18286693
	Chunga burmeisteri	0.3471905	0.34163687	-0.21079393	0.046399638
	Herpetotheres cachin...	0.49282703	0.13153149	-0.42115825	0.22466694
...					

# Data exploration with CLIP embeddings



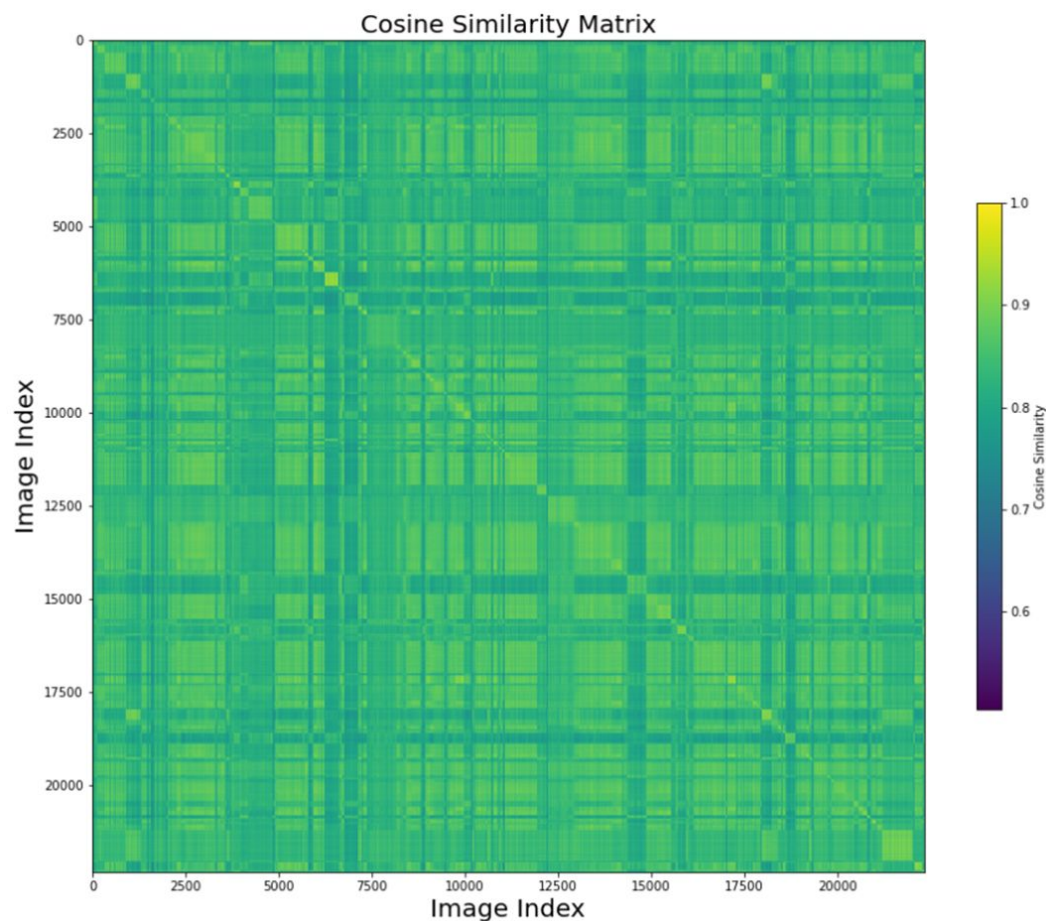
# Data exploration with CLIP embeddings



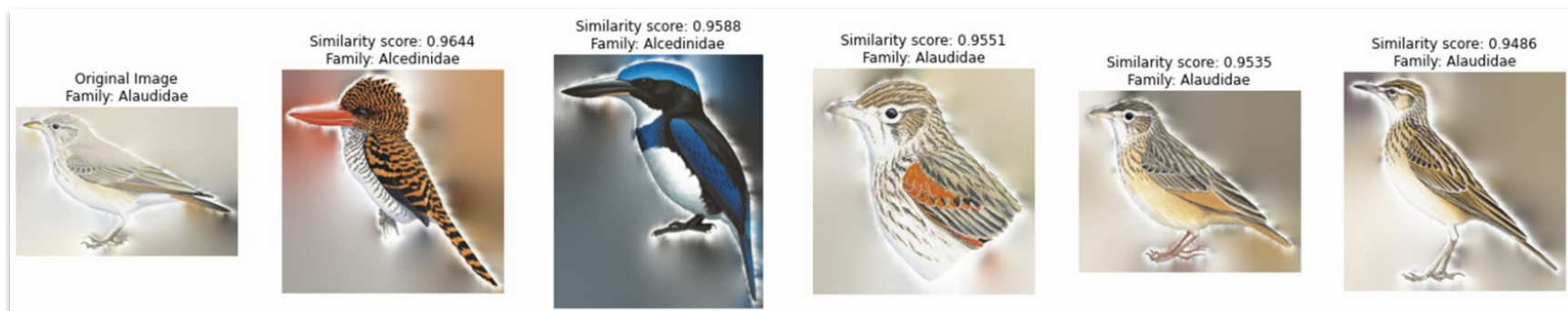
$$\text{similarity}(A, B) = \cos(\theta) = \frac{A \cdot B}{\|A\| \|B\|}$$



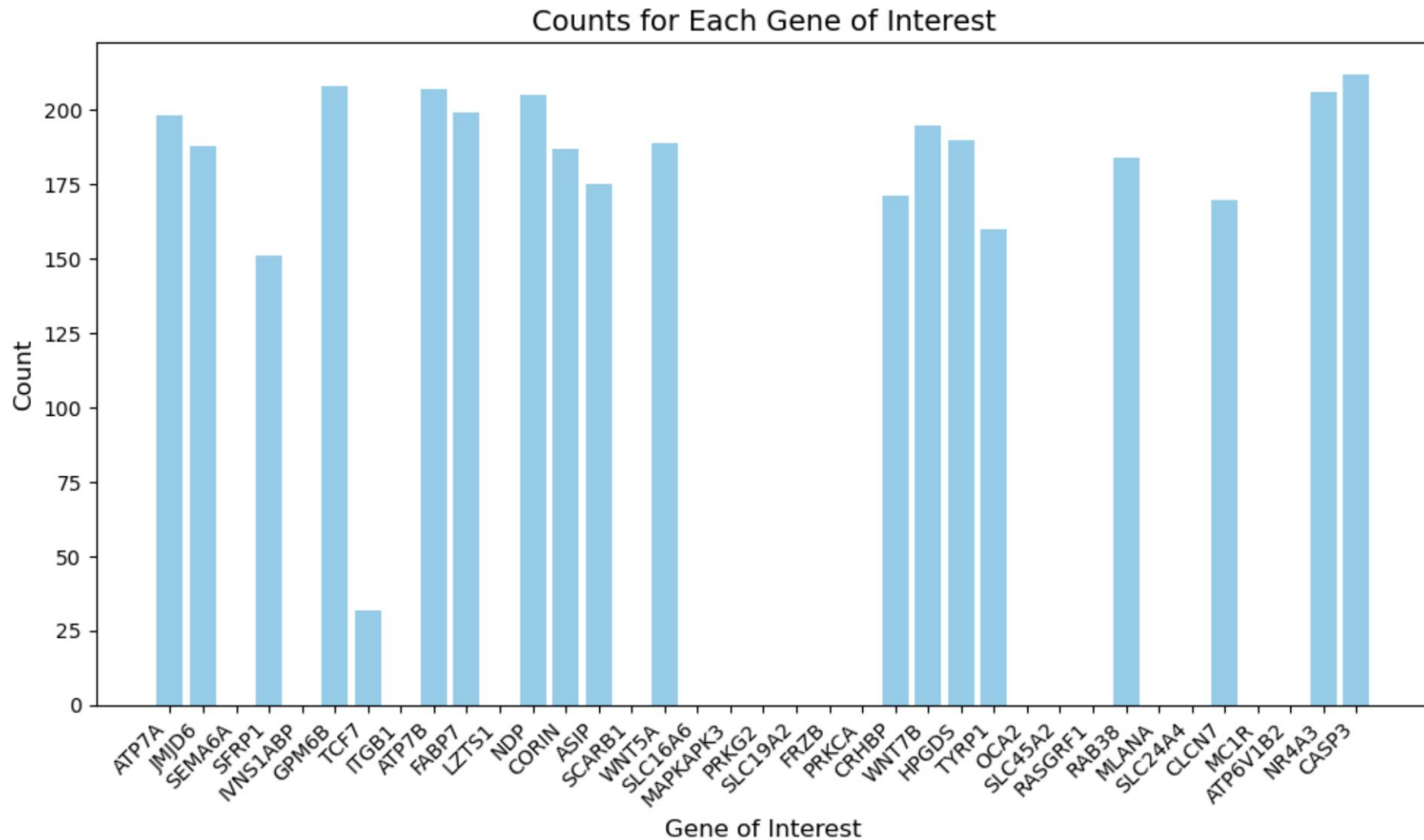
# Data exploration with CLIP embeddings



# Data exploration with CLIP embeddings



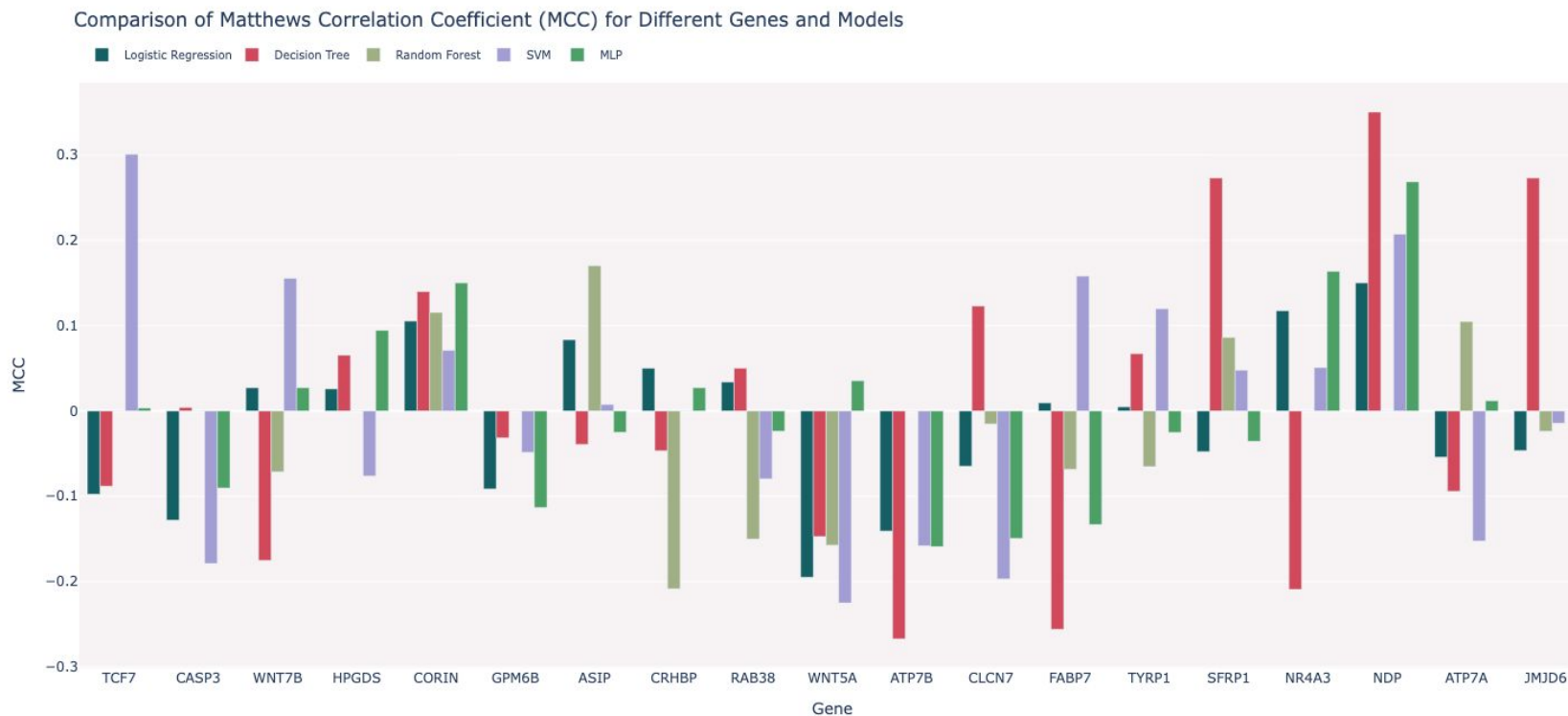
# Coloration genes



# Input for our classifiers

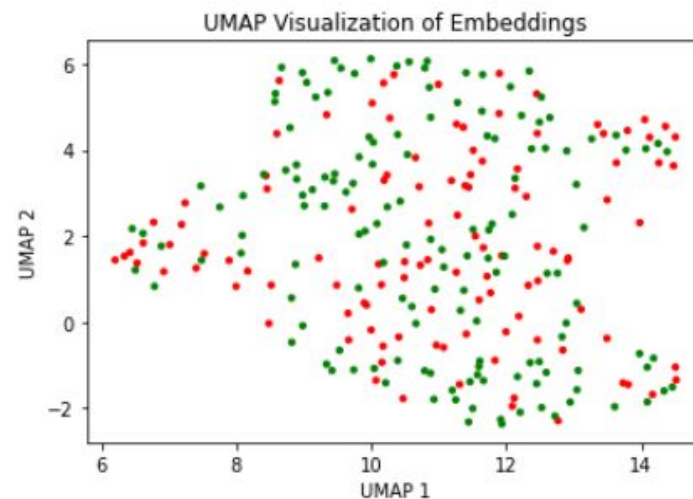
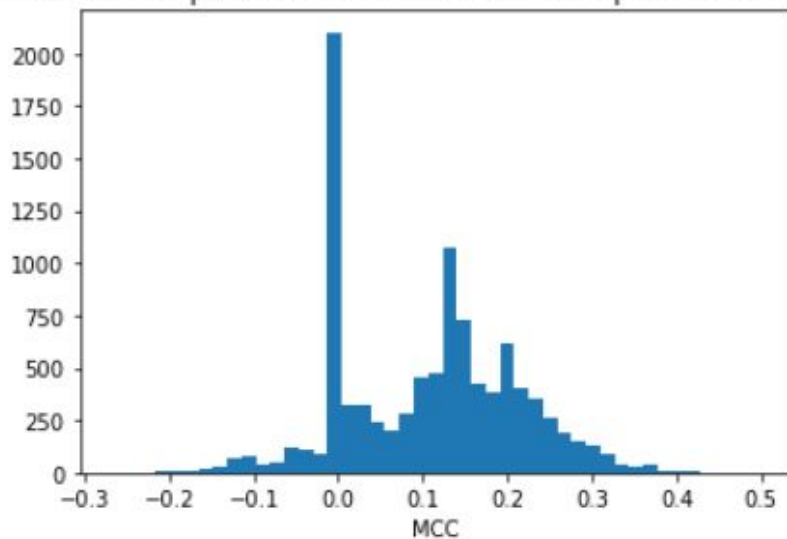
		All 512 CLIP embedding features				Y
All 258 species	Nothoprocta ornata	0.4072875	0.36193112	-0.31036815	0.13646944	1
	Smithornis capensis	0.3131959	0.20072602	-0.43035	0.18086748	0
	Formicarius rufipectus	0.37125307	-0.050516717	-0.52256465	0.25555477	1
	Sylvia atricapilla	0.21032035	0.24915919	-0.6387599	0.5333118	0
	Lanius ludovicianus	0.6003466	0.18232825	0.022367803	0.04154638	1
	Amazona guildingii	0.33798376	0.14806776	-0.47427034	-0.19996244	1
	Probosciger aterrimus	0.36870858	-0.056442954	-0.33242458	-0.04306676	0
	Eolophus roseicapilla	0.40817946	0.23949404	-0.13612387	0.18286693	0
	Chunga burmeisteri	0.3471905	0.34163687	-0.21079393	0.046399638	1
	Herpetotheres cachin...	0.49282703	0.13153149	-0.42115825	0.22466694	...

# Classical ML models

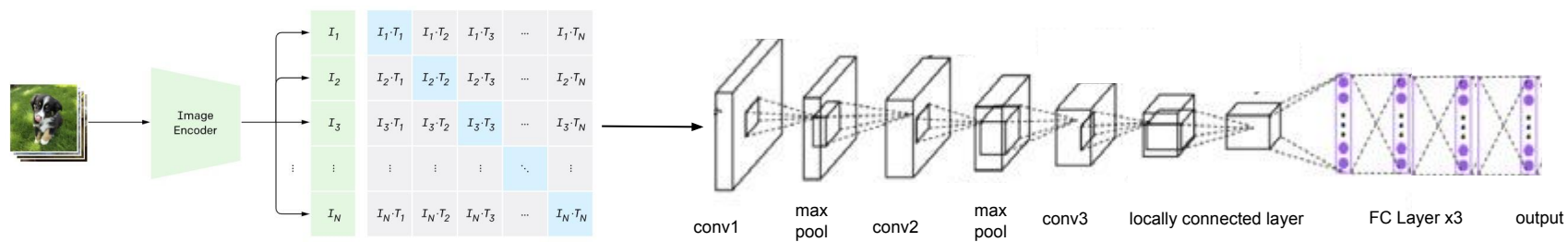


# Bootstrapping

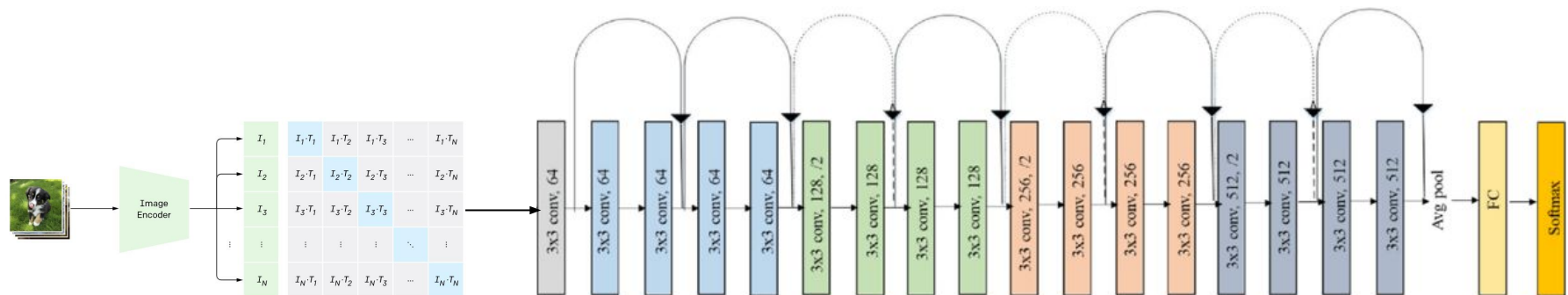
MCC for the AKAP5 protein for 10000 different Datasplits. Mean MCC = 0.107



# Gene2Bird

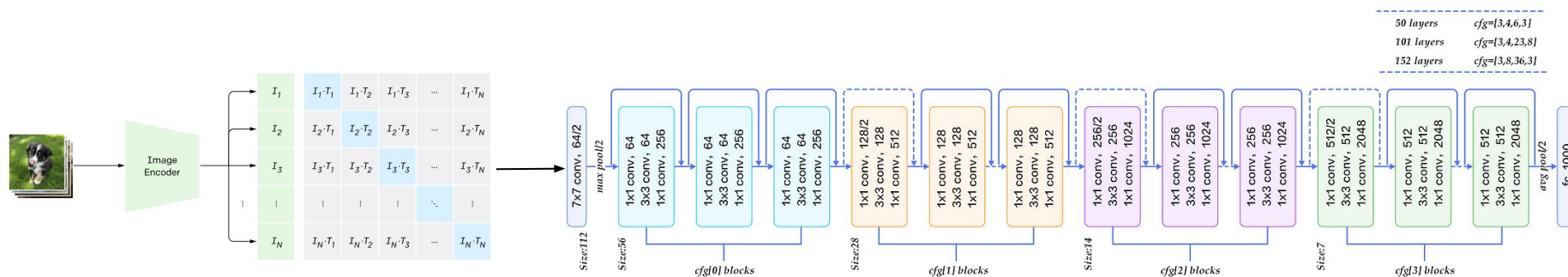


# ResNet18

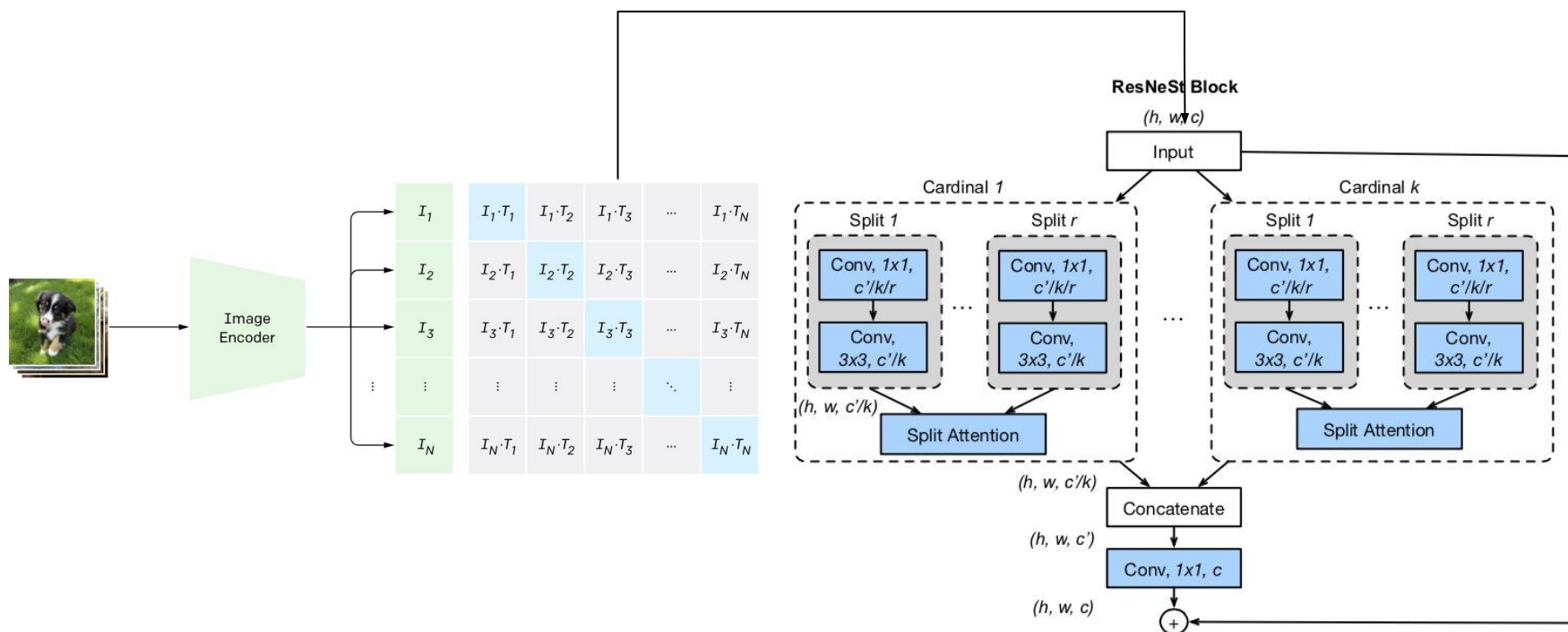




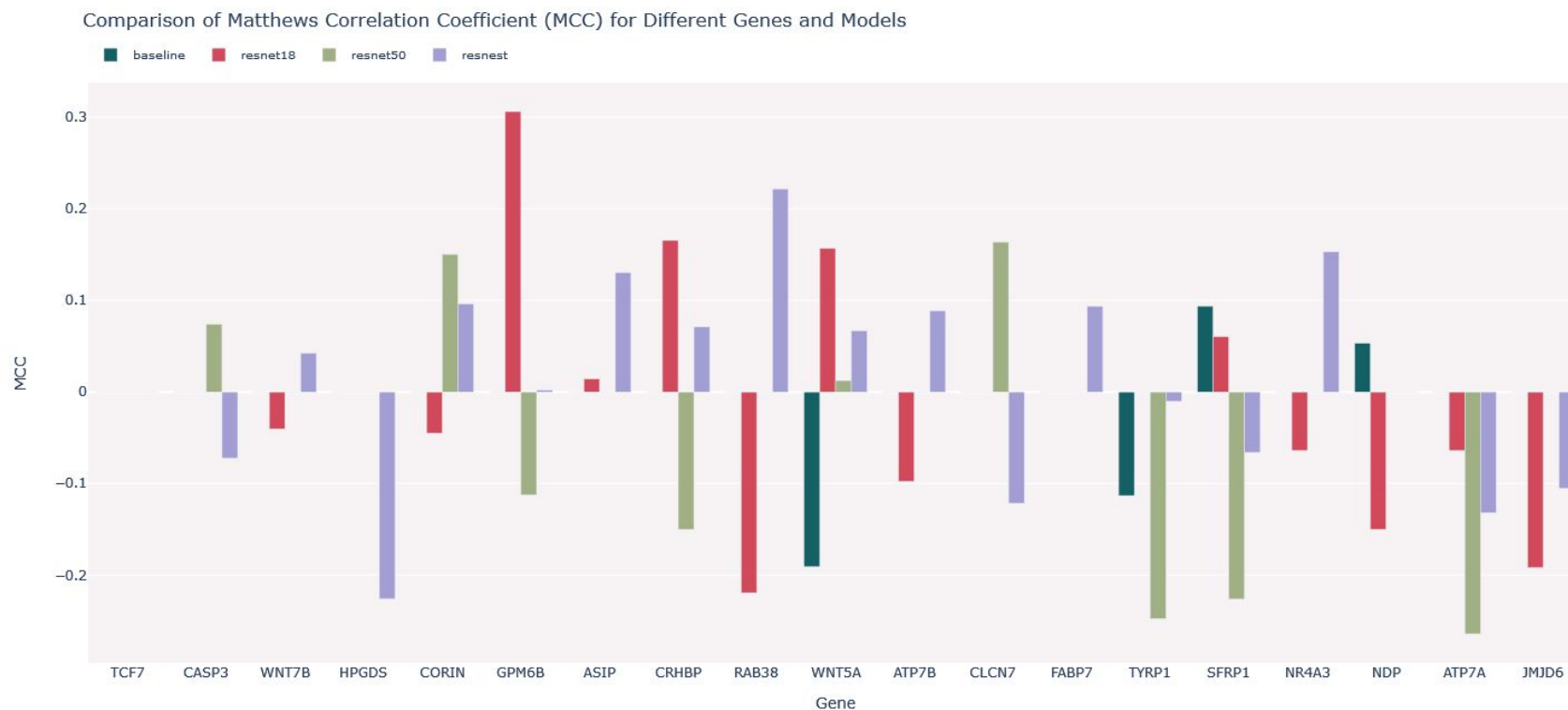
# ResNet50



# ResNeSt



# Results of Deep Learning Models



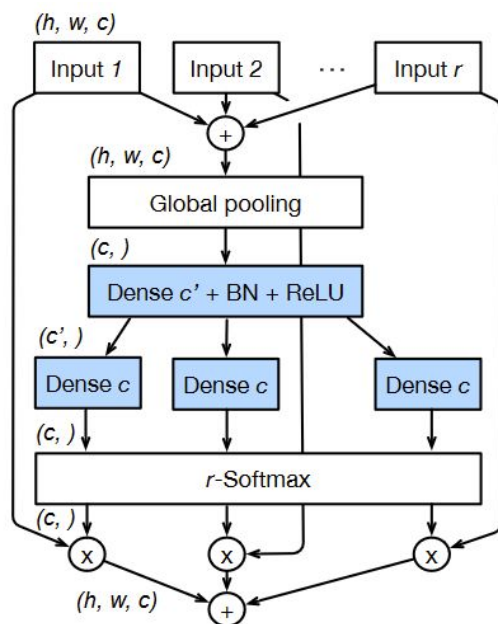
# Potential Issues

- CLIP might not have captured all the necessary information
- Poor genbank annotation
- Presence or absence of a gene may not provide sufficient information

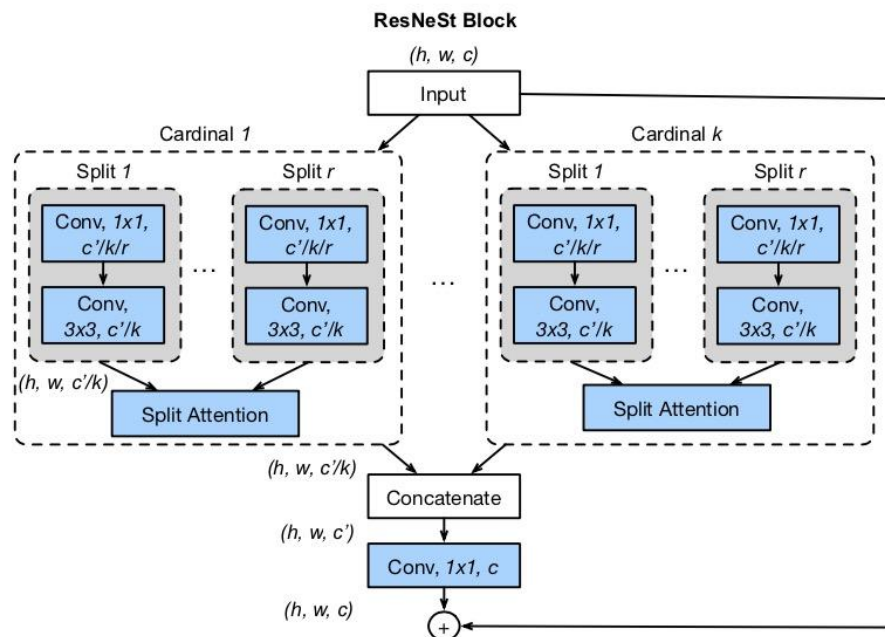
**Thank you for your time and attention!**

**Any questions?**

# ResNeSt - Details



Split-Attention within a cardinal group



ResNeSt block