

CELL-E 2: Translating Proteins to Pictures and Back with a Bidirectional Text-to-Image Transformer

Emaad Khwaja

UC Berkeley - UCSF Bioengineering Graduate Program

UC Berkeley EECS

Goal

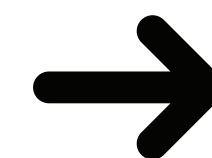
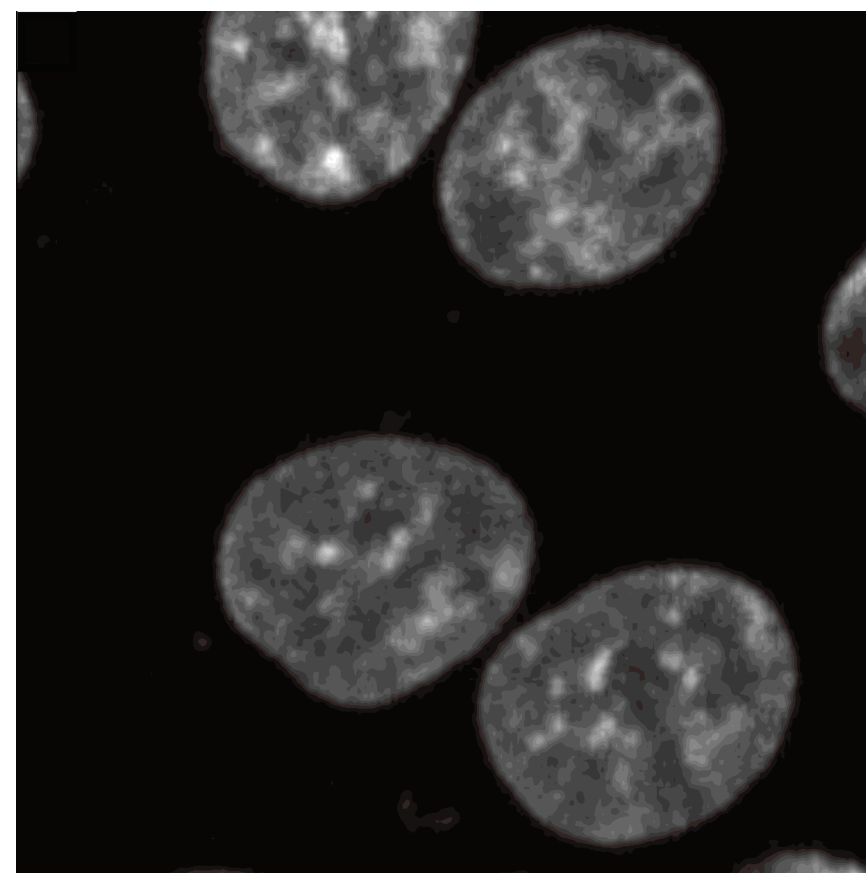
Predict protein localization *in silico*.

Protein Sequence

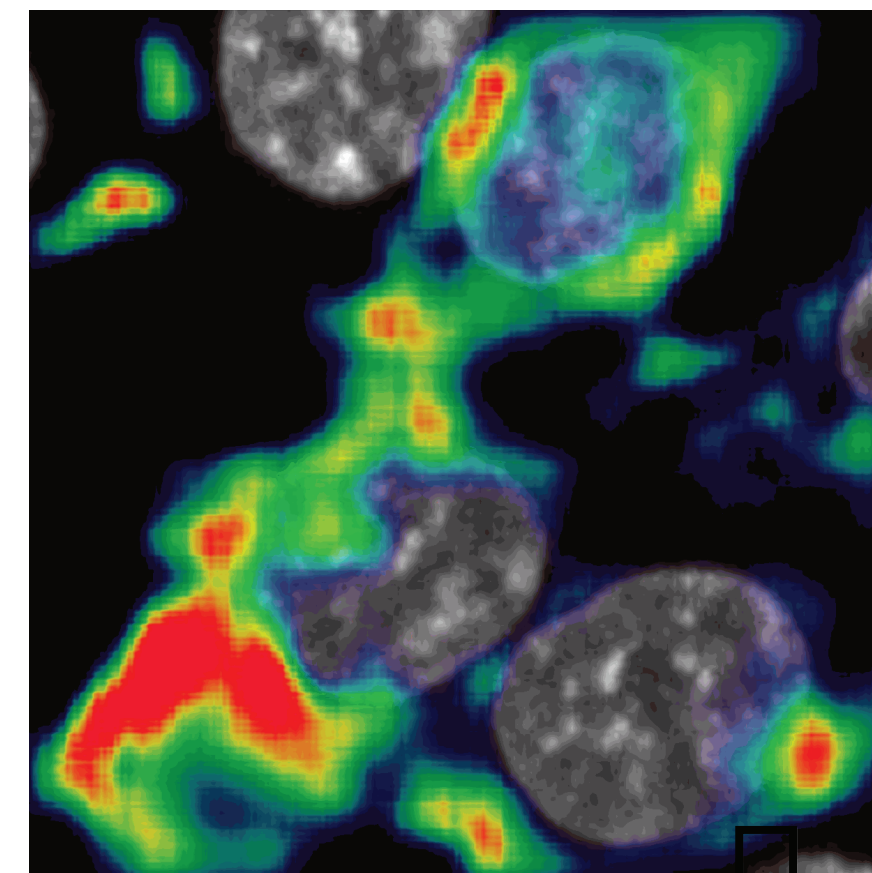
MVCFRLFPVPGSGLVLV
CLVLGAVRSYALELNL
DSENATCLYAKWQMNF
TVRYETTNKTYKTVTI
SDHGTVTYNGSICGDD
QNGPKIAVQFGPGFSWI
ANFTKAASTYSIDSVSF
SYNTGDNTTFPDAEDK
GILTVDELLAIRIPLND...

+

Nucleus Image



Predicted Localization



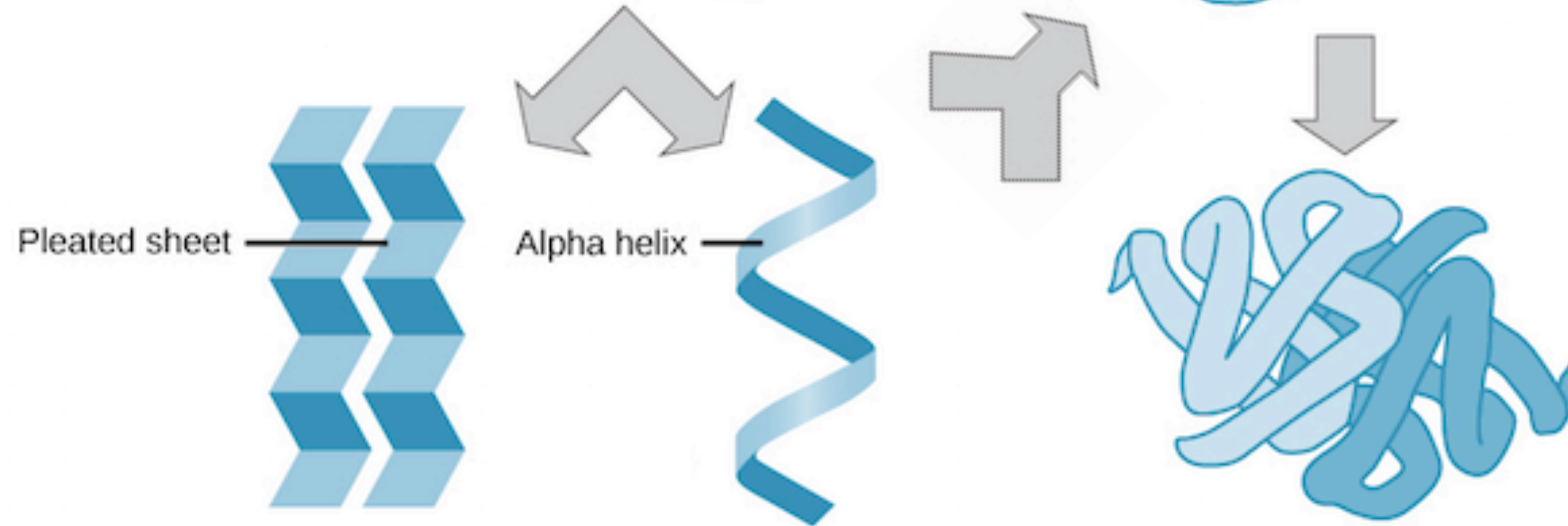
Why?

Proteins dictate biological function.

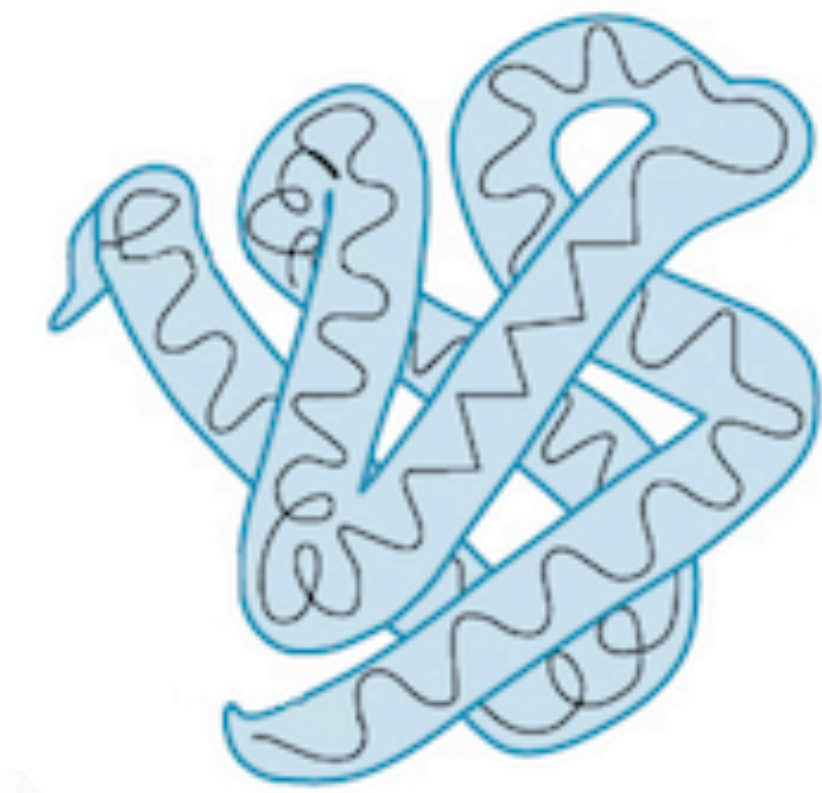
Primary protein structure
sequence of a chain of amino acids



Secondary protein structure
hydrogen bonding of the peptide backbone causes the amino acids to fold into a repeating pattern



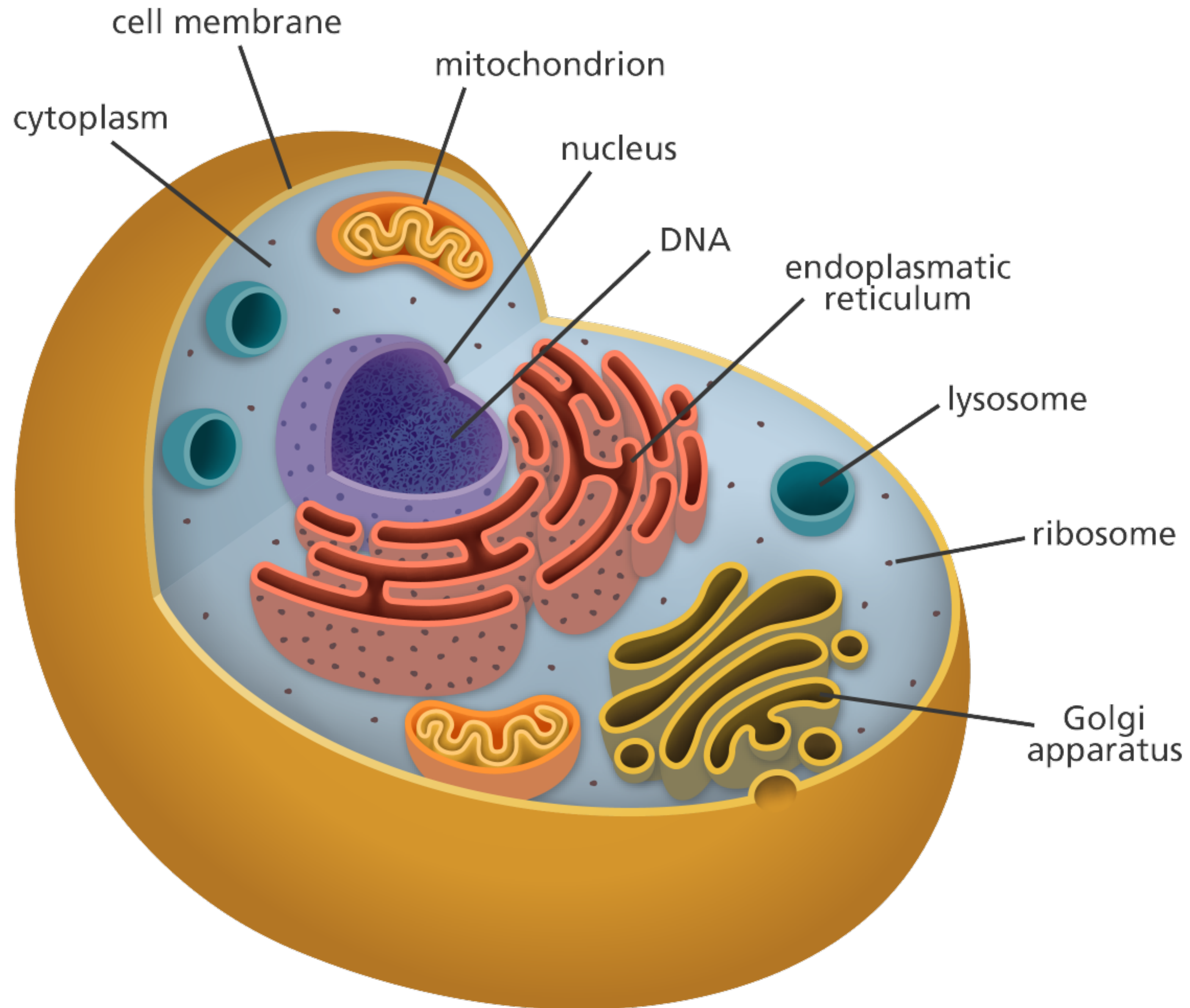
Tertiary protein structure
three-dimensional folding pattern of a protein due to side chain interactions



Quaternary protein structure
protein consisting of more than one amino acid chain

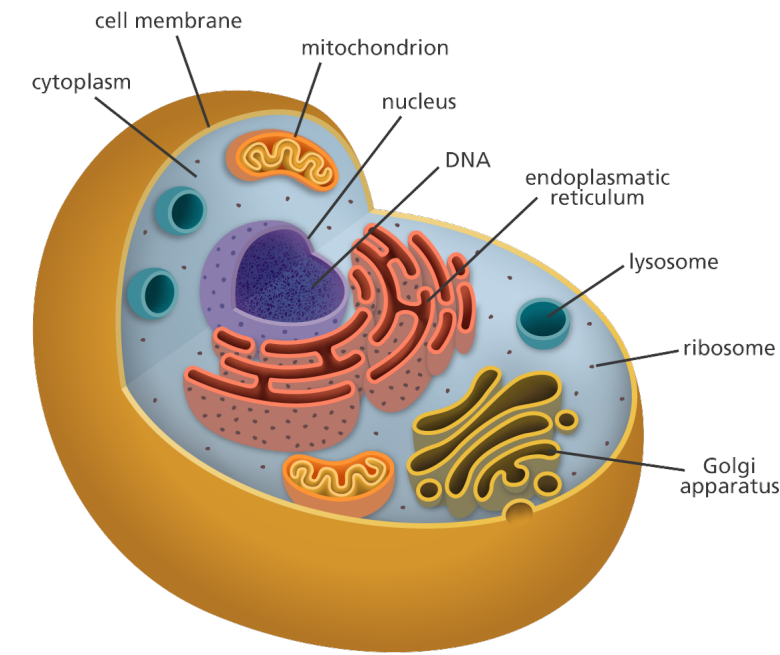


Predicting localization allows us
to gain fundamental
understanding of biology.



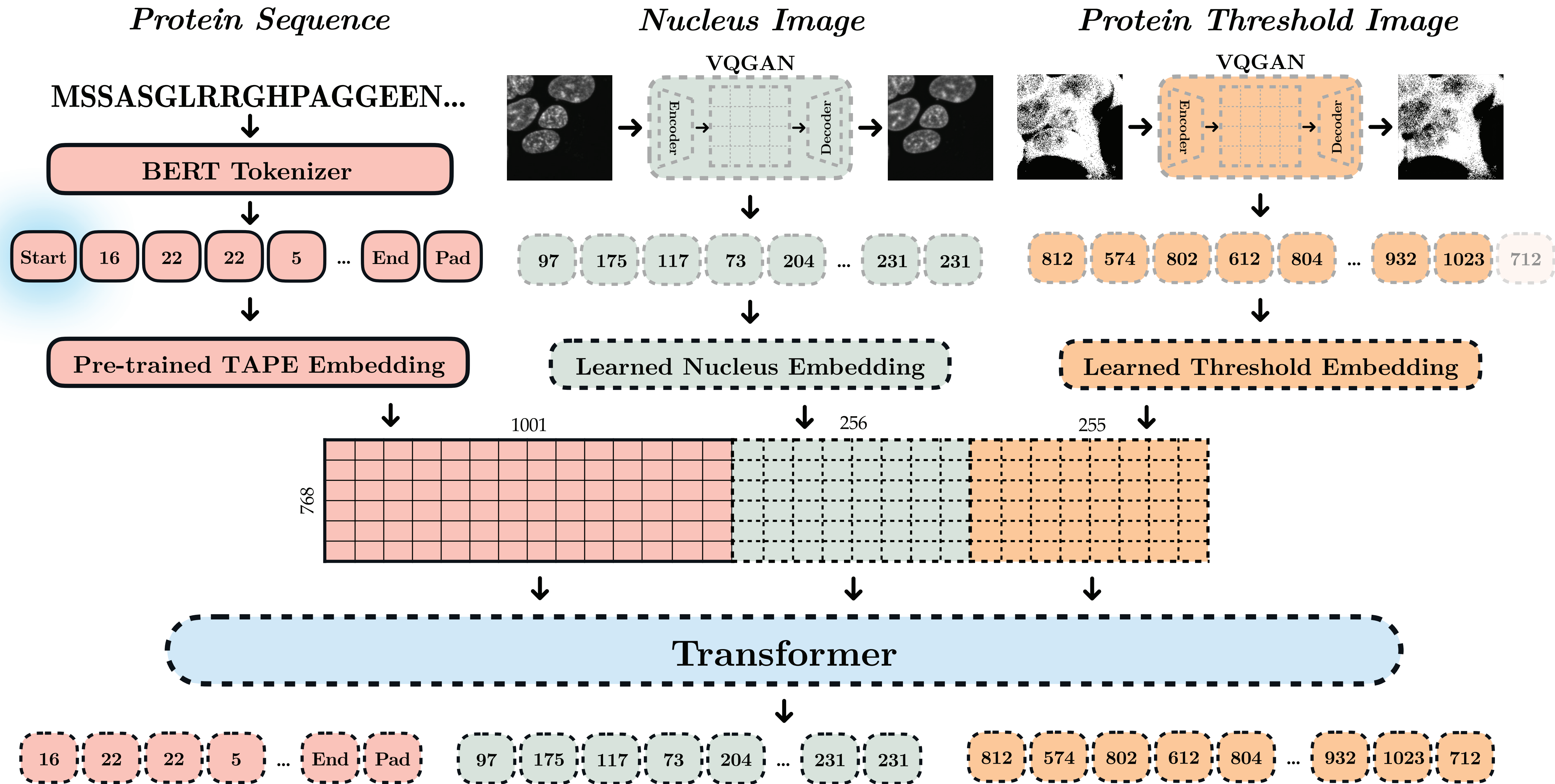
Source: Genome Research Limited

**Correct localization is necessary
for function.**



**Mislocalization can result in
dysfunction and disease.**

CELL-E



DNA Topoisomerase I
Relative Attention Weights

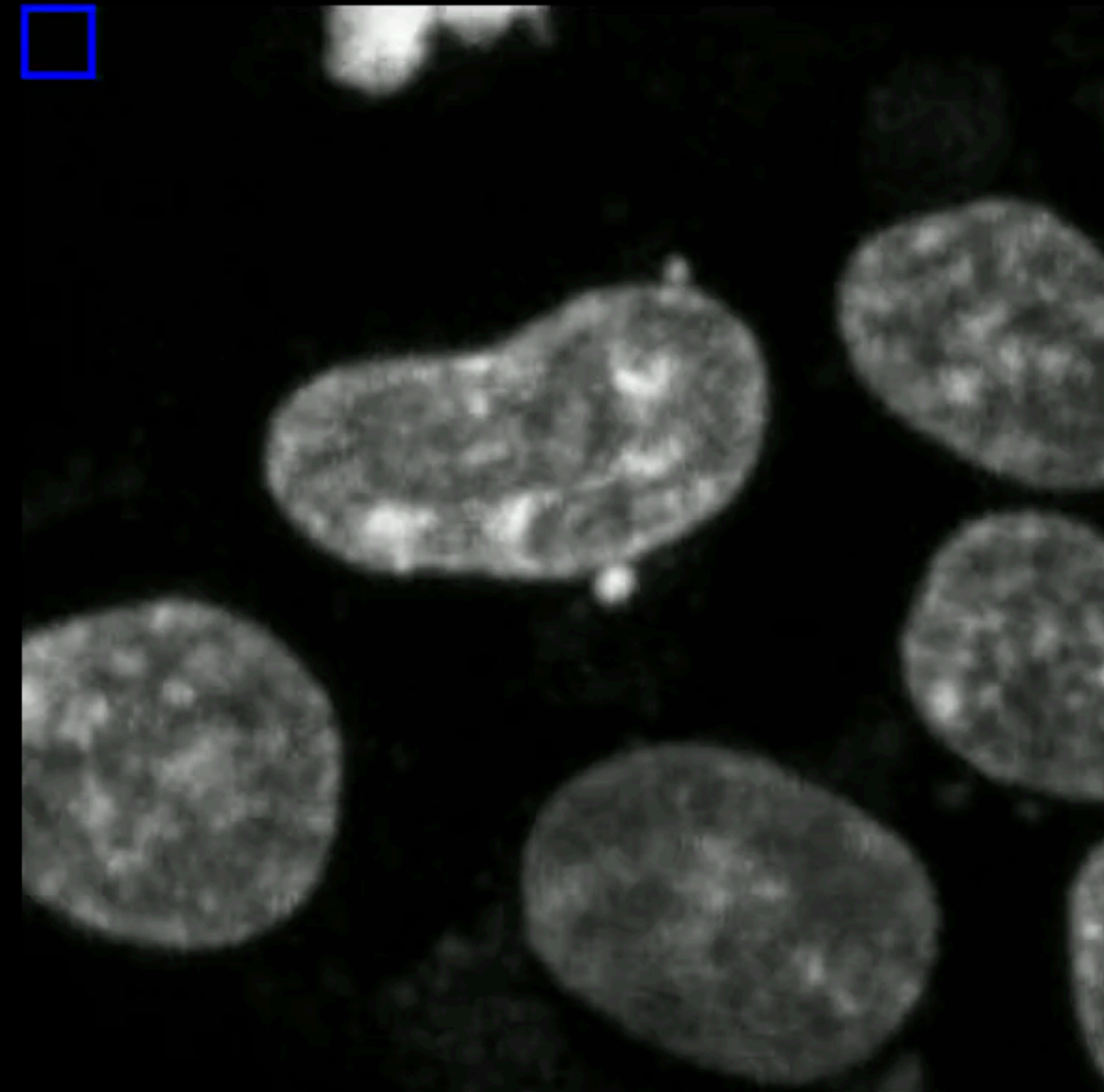
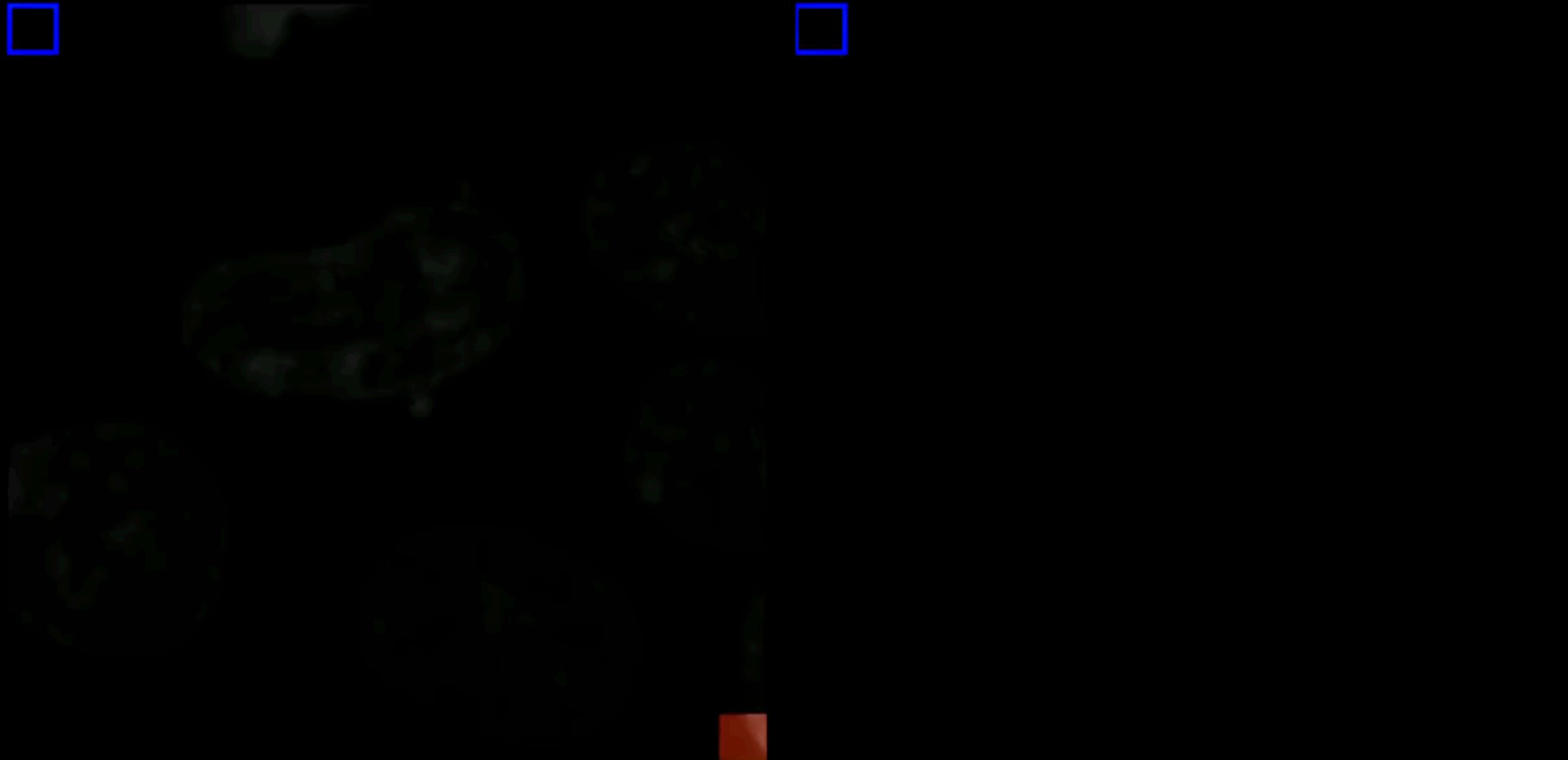
Sequence

Nucleus Image

Predicted Threshold Image

Predicted PDF

```
MSGDHLHND SQ IEAD FRLNDS HKHKDKH  
KDREHRHKEHKK EKDR EKSKHSNSEHKD  
SEK KHKEKEKT KHK DGSSEKHKDKHKDR  
DKEKRKEEKVRASGDAKIKKEKENGFS  
PPQIKDEPEDDGYFVPPKE DIKPLKRPR  
DEDDADYKPKKIKTEDTKKEKKRKL EEE  
EDGK LKKPKNKDKDKKVPEPDNKKKKPK  
KEEEQKWKWEEERYPEGIKWKFLEHKG  
PVFAPPYEPLPENVKFYD GKVMKLSPK  
AEEVATFFAKMLDHEYTTKEIFRKNFFK  
DWRKEMTNEEKNIITNLSKCDFTQMSQY  
FKAQTEARKQMSKEEKLKIKEENEKLLK  
EYGF CIMDNHKERIANFKIEPPGLFRGR  
GNHPKMGMLKRRIMPEDI IINCSKDAKV  
PSPPPGHKWKEVRHDNKVTWLVS WTENI  
QGSIKYIMLNPS SR IKGEKDWQKYETAR  
RLKKCVDKIRNQYREDWKS KEMKVRQRA  
VALYFIDK LALRAGNEKEEGETADTVGC  
CSLRVEHINLHP ELDGQEYVVEFD FLGK  
DSIRY YNKVPVEKRVFKNLQLFMENKQP  
EDDLFDRLNTGILNKHLQDLMEGLTAKV  
FRTYNASITLQQQLKELTAPDENIPAKI  
LSYNRANRAVA ILCNHQRAPPKTFEKSM  
MNLQTKIDAKKEQLADARRDLKS AKADA  
KVMKDAKTKKVVESKKA VQRLEEQLMK  
LEVQATDREENKQIALGTSKLN YLDPRI  
TVAWCKKWGVPIEKIYNKTQREKFAWA I  
DMADEDE
```



CELL-E 2

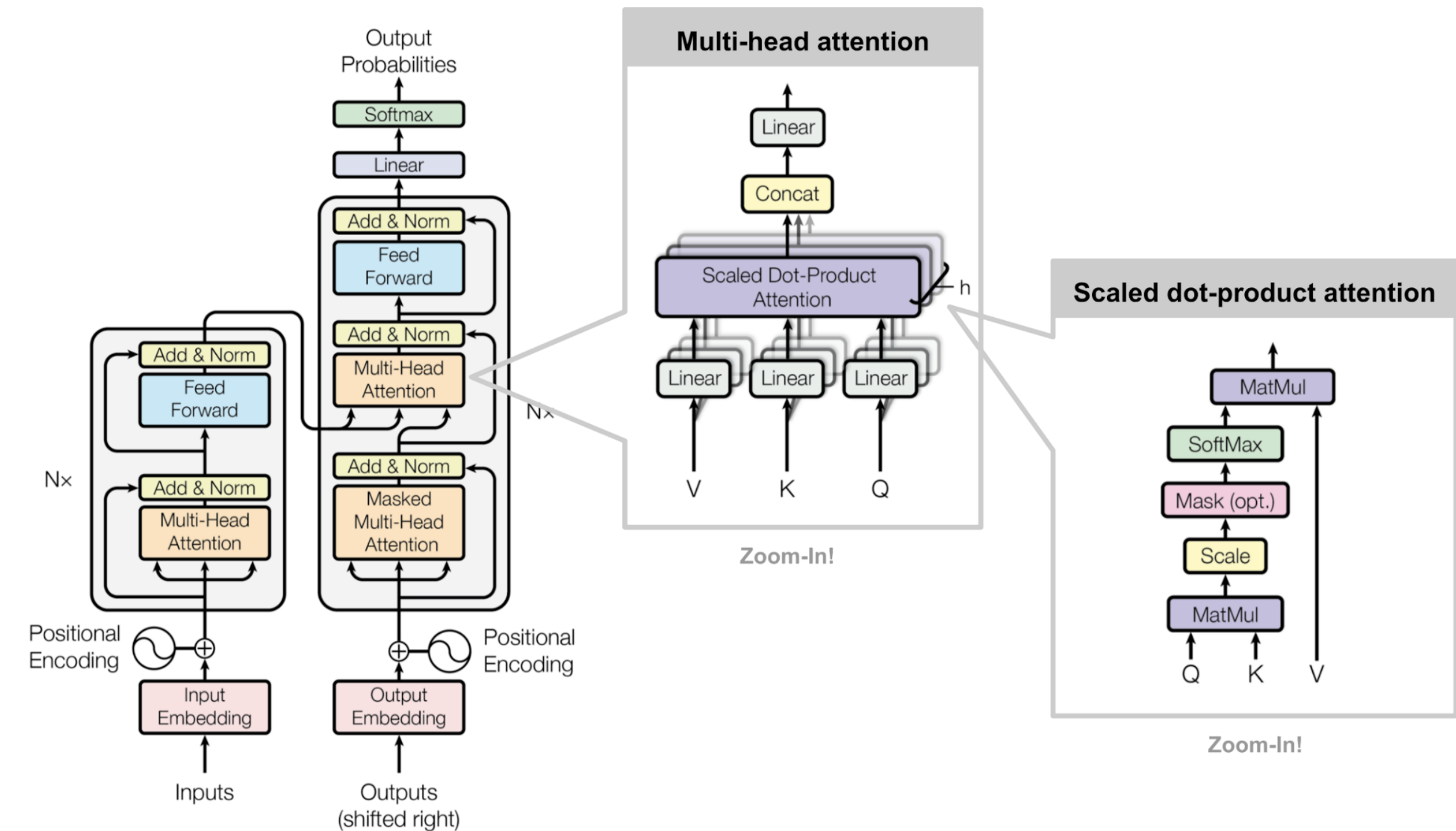
Drawbacks

- **Small dataset**
 - **Not many proteins in OpenCell**
- Unidirectional
 - Can do text-to-image but cant tell you what text should be there with a localization pattern
- Autoregressive generation is slow

THE HUMAN PROTEIN ATLAS

Drawbacks

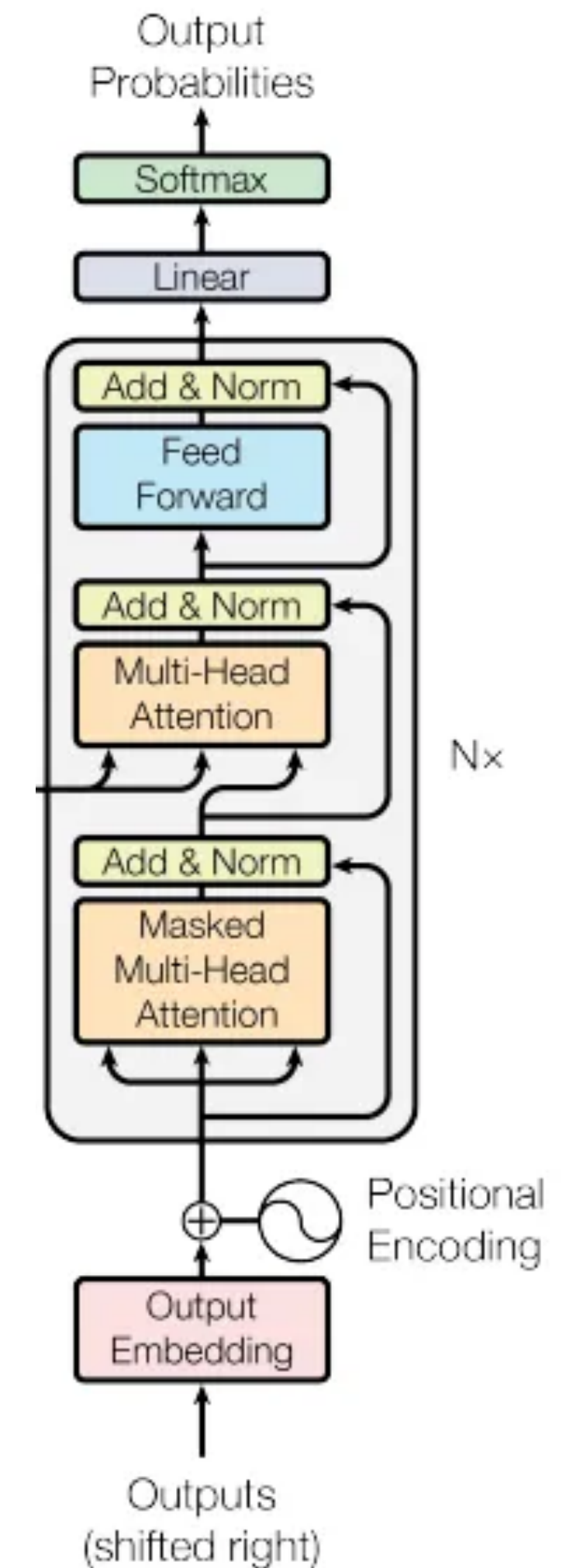
- Small dataset
 - Not many proteins in OpenCell
- **Unidirectional**
 - **Can do text-to-image but cant tell you what text should be there with a localization pattern**
- Autoregressive generation is slow



Drawbacks

- Small dataset
 - Not many proteins in OpenCell
- **Unidirectional**
 - **Can do text-to-image but cant tell you what text should be there with a localization pattern**
- Autoregressive generation is slow

Decoder-Only

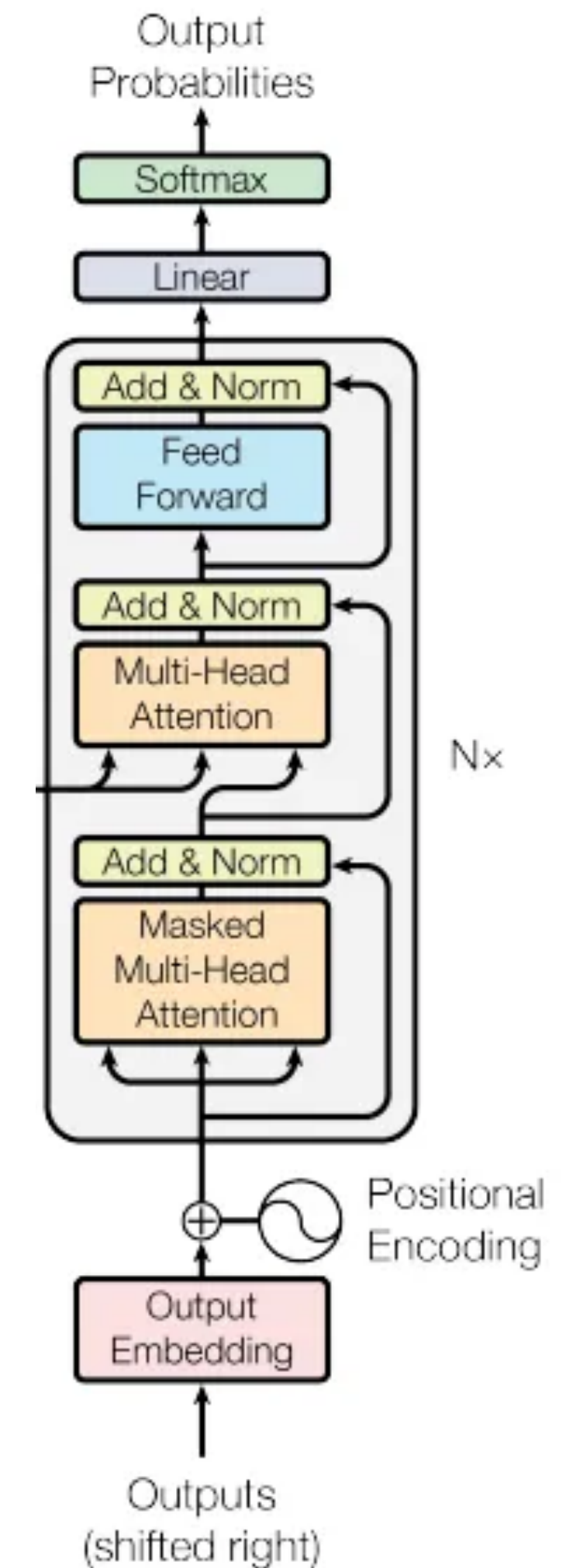


Drawbacks

- Small dataset
 - Not many proteins in OpenCell
- **Unidirectional**
 - **Can do text-to-image but cant tell you what text should be there with a localization pattern**
- Autoregressive generation is slow

Single Token Output

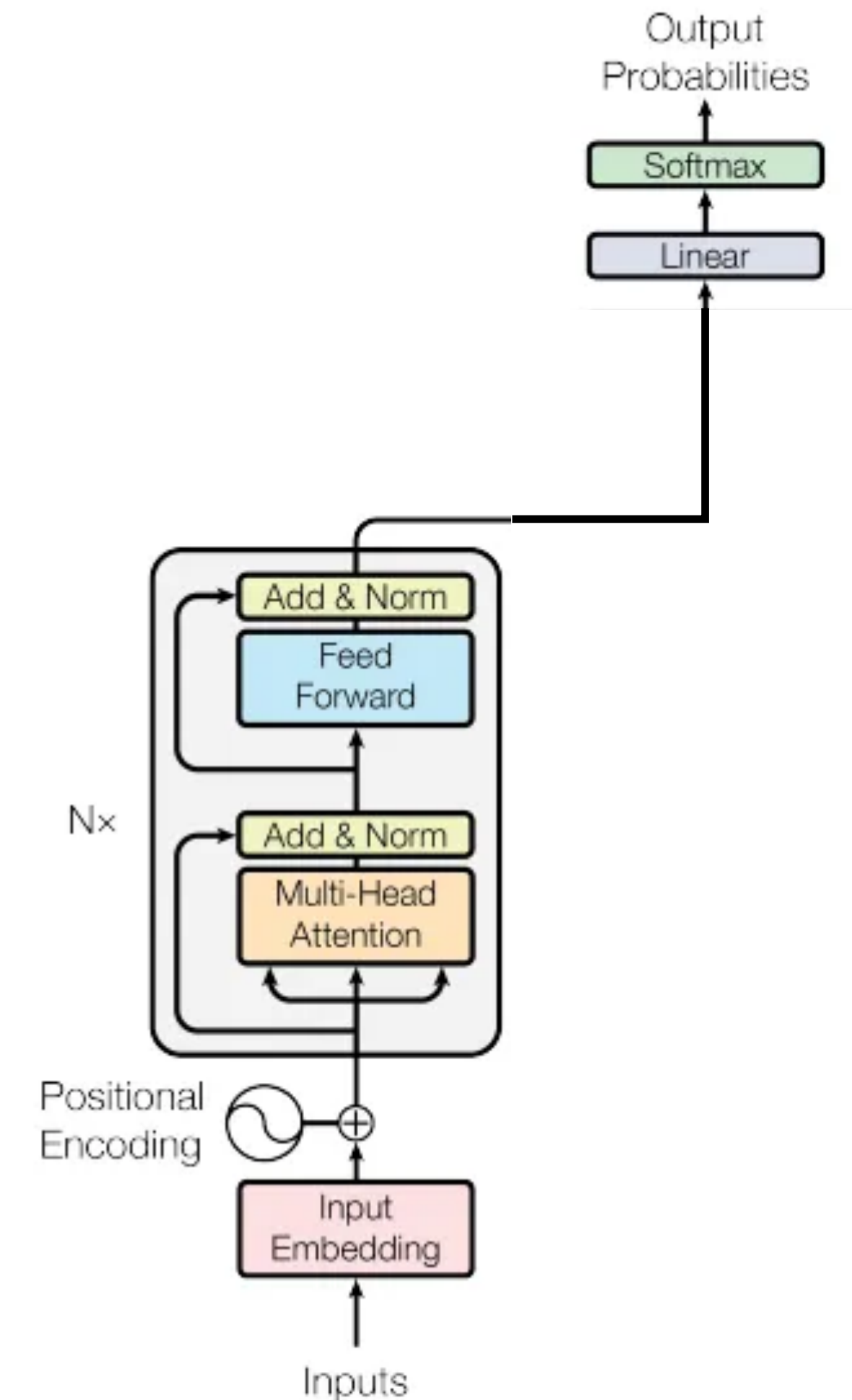
Decoder-Only



Drawbacks

- Small dataset
 - Not many proteins in OpenCell
- **Unidirectional**
 - **Can do text-to-image but cant tell you what text should be there with a localization pattern**
- Autoregressive generation is slow

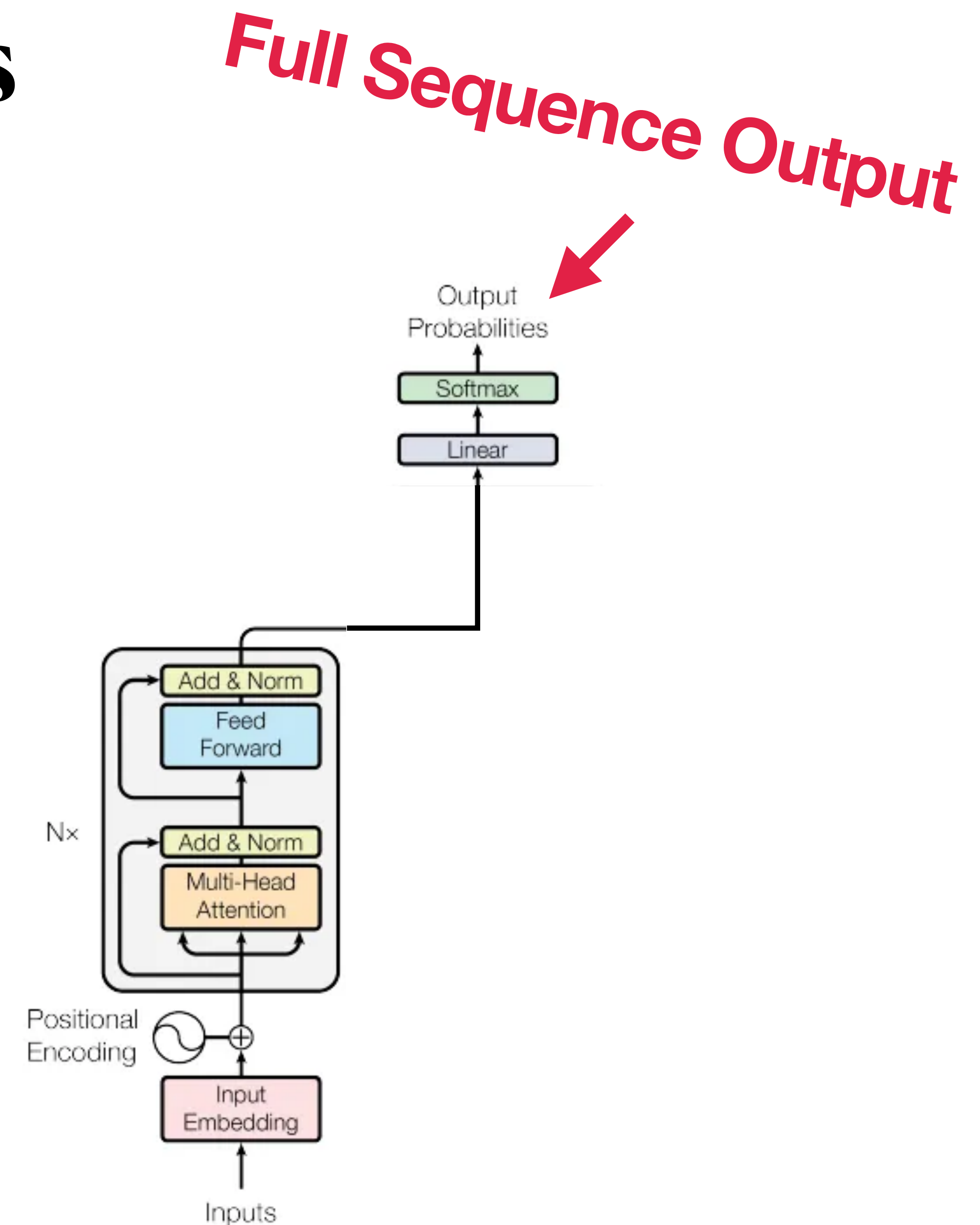
Encoder-Only



Drawbacks

- Small dataset
 - Not many proteins in OpenCell
- **Unidirectional**
 - **Can do text-to-image but cant tell you what text should be there with a localization pattern**
- Autoregressive generation is slow

Encoder-Only

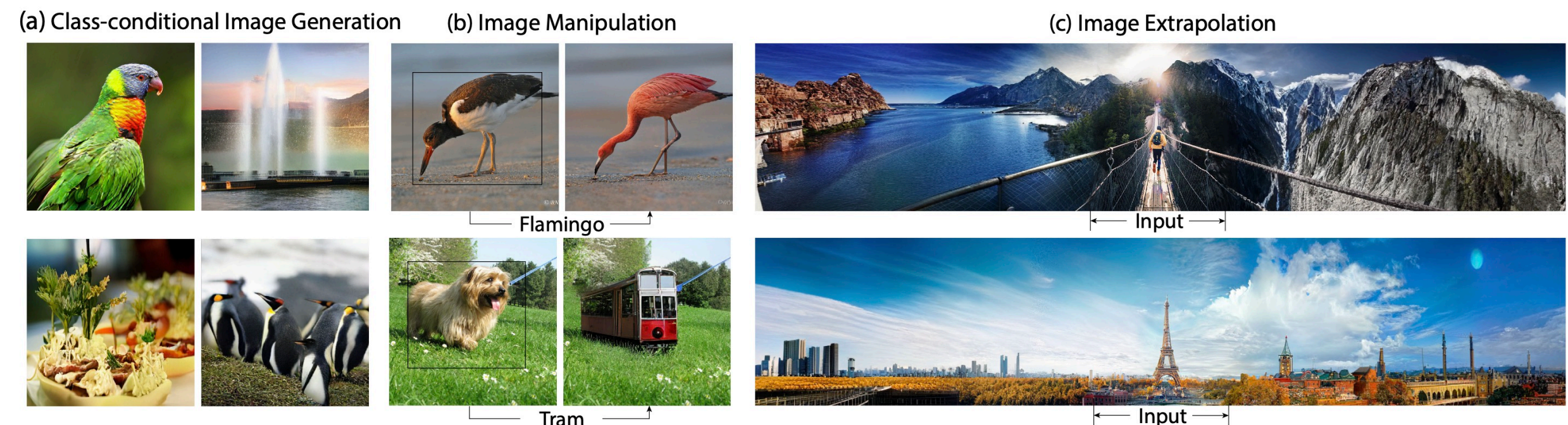


Drawbacks

- Small dataset
 - Not many proteins in OpenCell
- Unidirectional
 - Can do text-to-image but cant tell you what text should be there with a localization pattern
- **Autoregressive generation is slow**

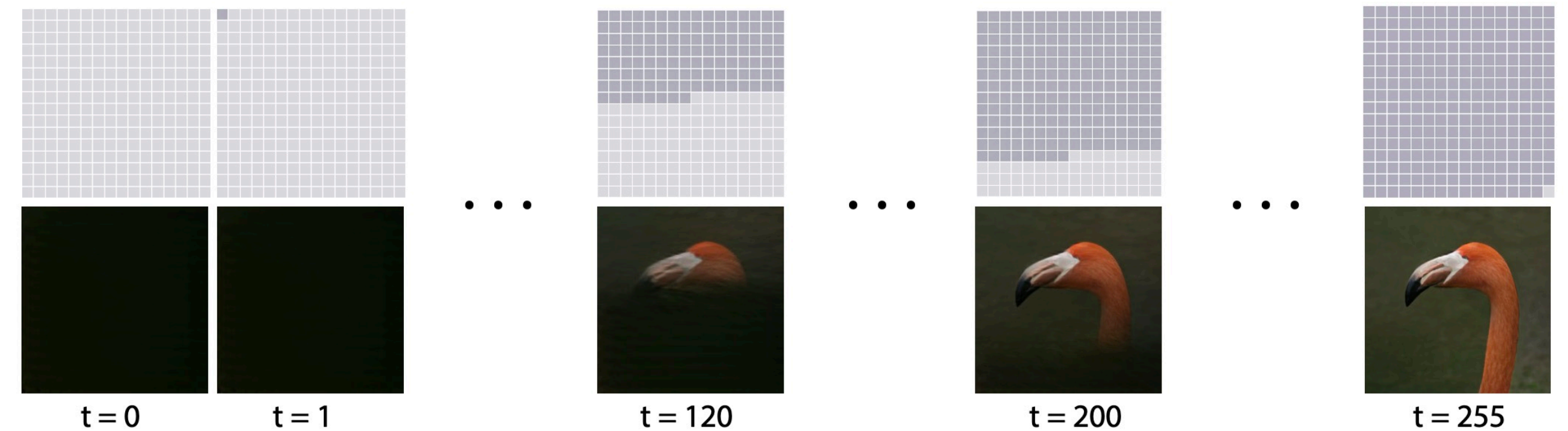
MaskGIT: Masked Generative Image Transformer

Huiwen Chang Han Zhang Lu Jiang Ce Liu* William T. Freeman
Google Research



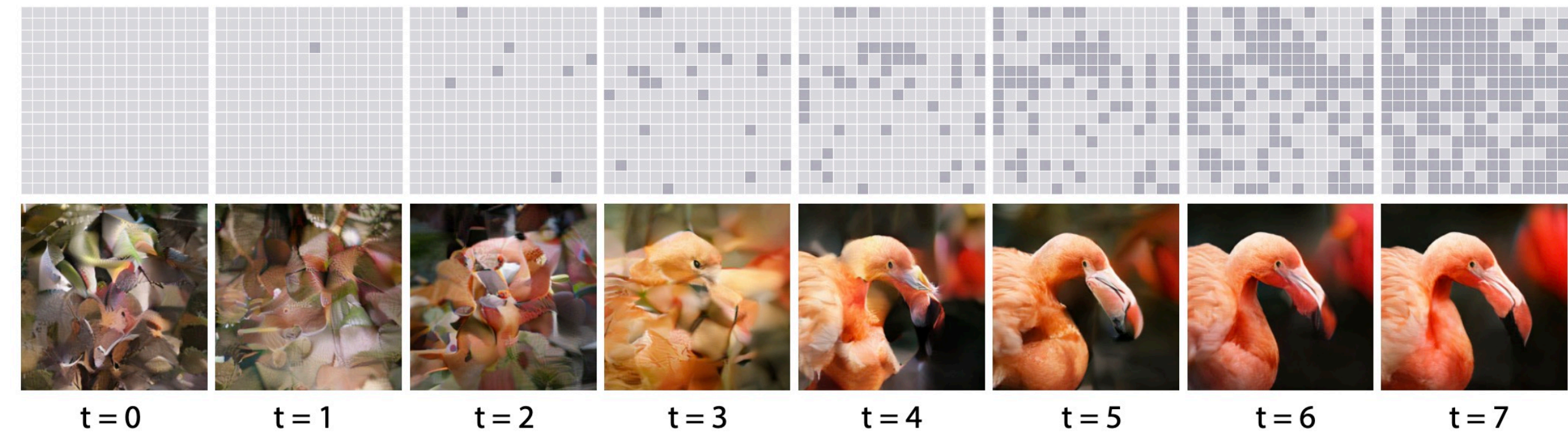
Drawbacks

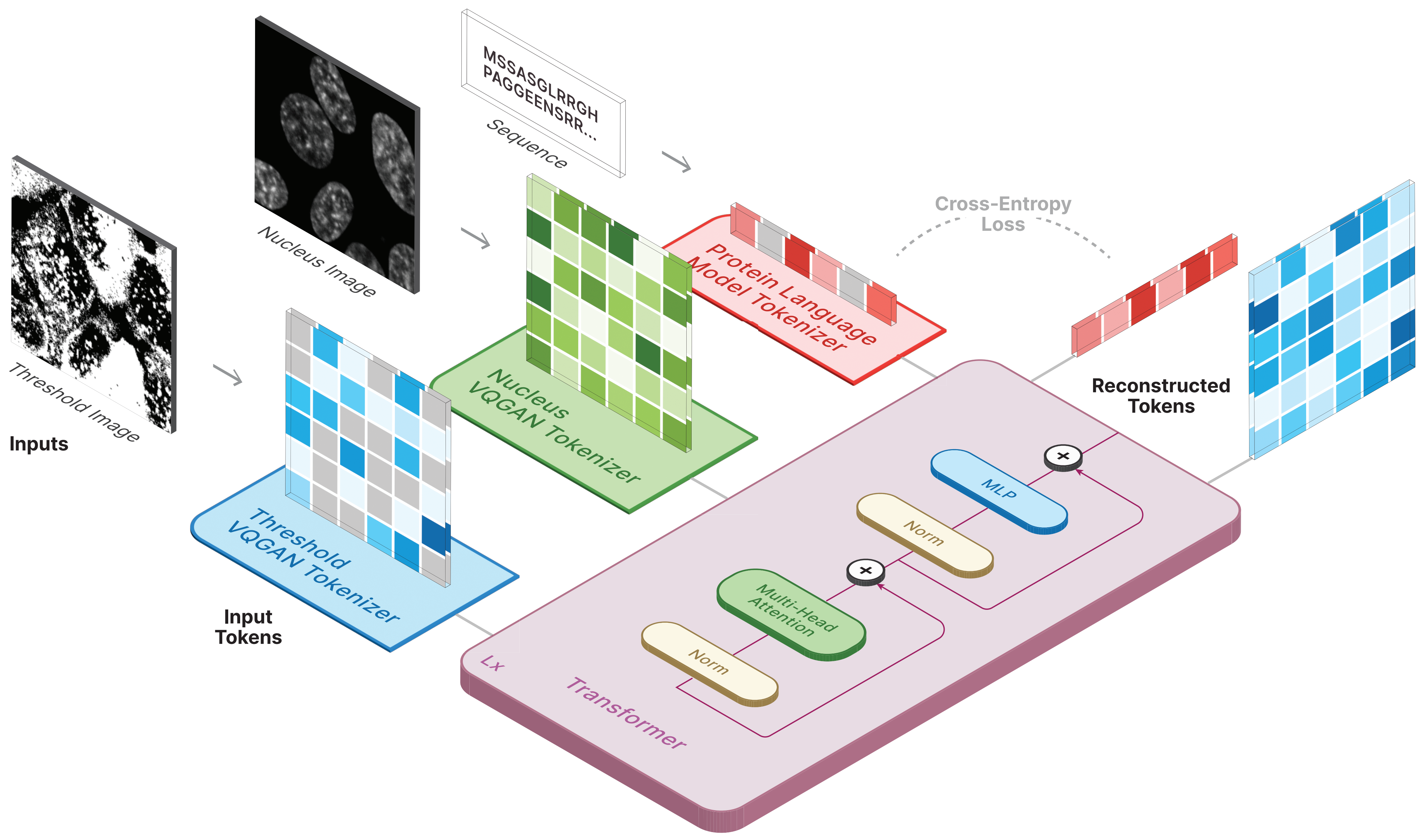
- Small dataset
 - Not many proteins in OpenCell
- Unidirectional
 - Can do text-to-image but cant tell you what text should be there with a localization pattern
- **Autoregressive generation is slow**

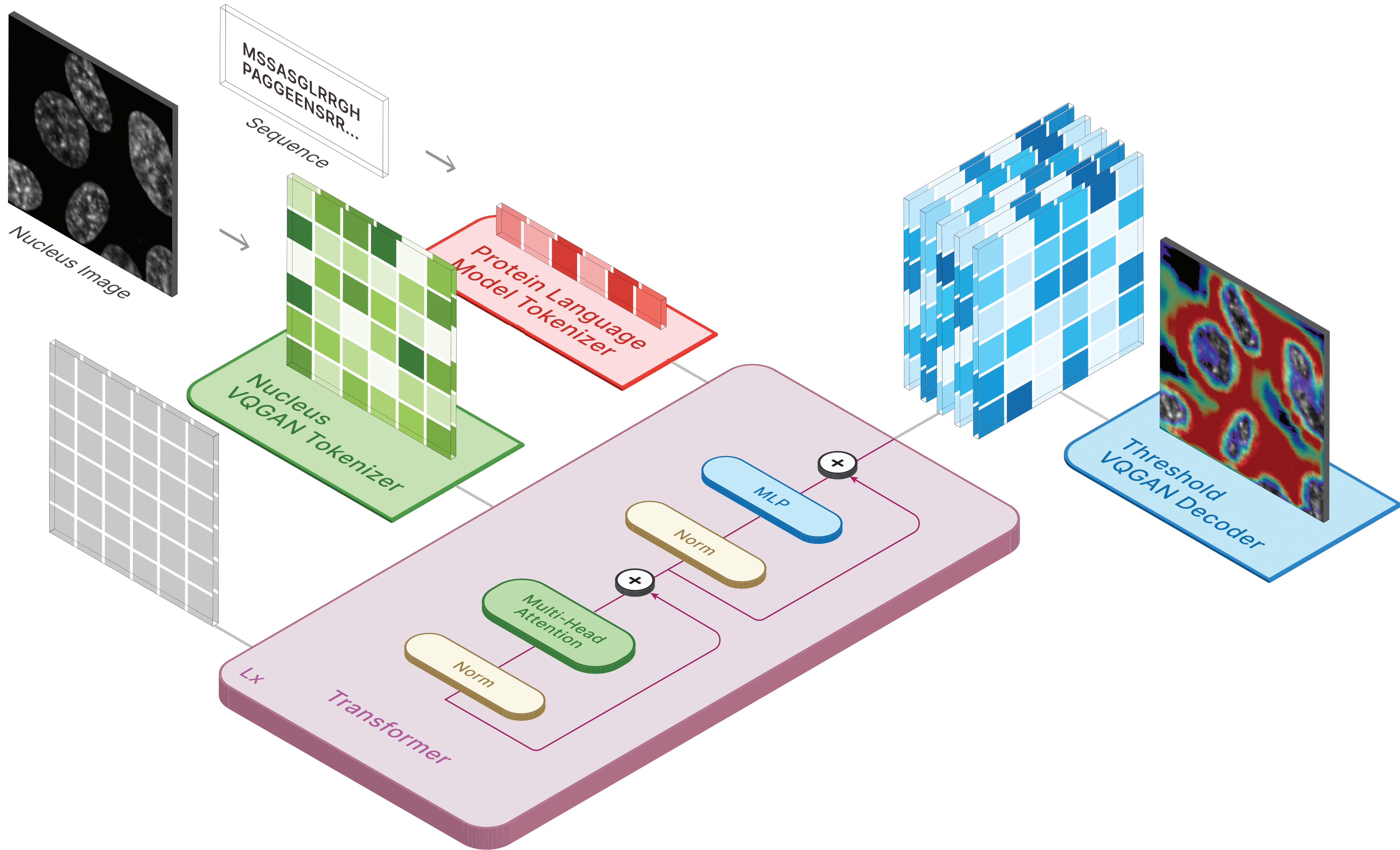


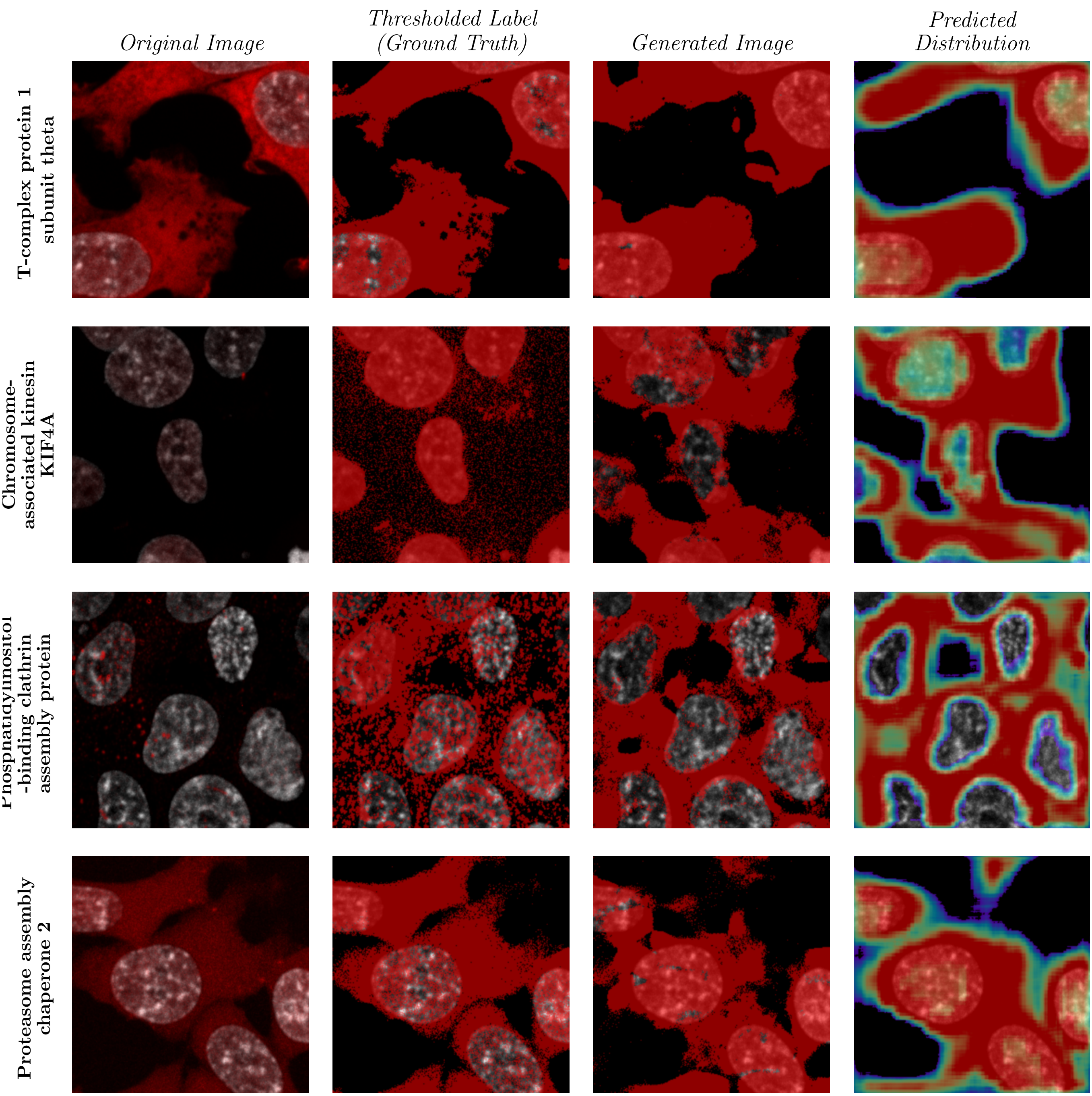
Drawbacks

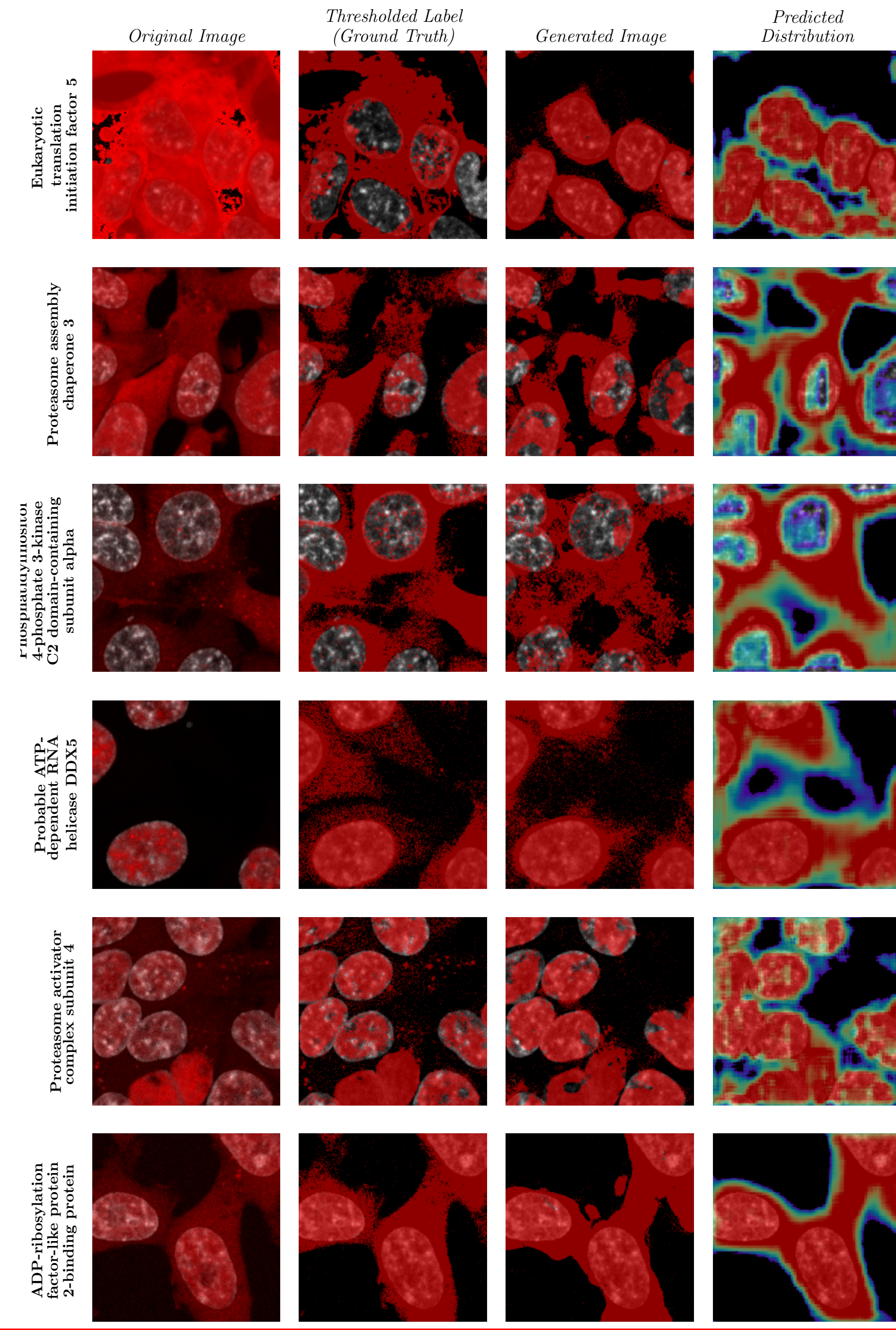
- Small dataset
 - Not many proteins in OpenCell
- Unidirectional
 - Can do text-to-image but cant tell you what text should be there with a localization pattern
- **Autoregressive generation is slow**





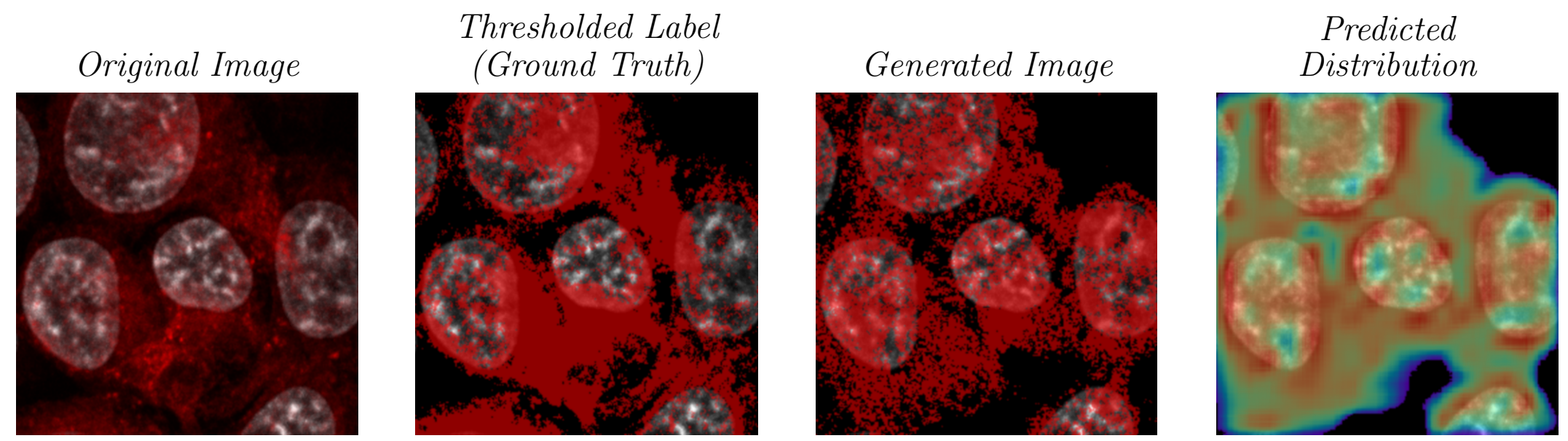




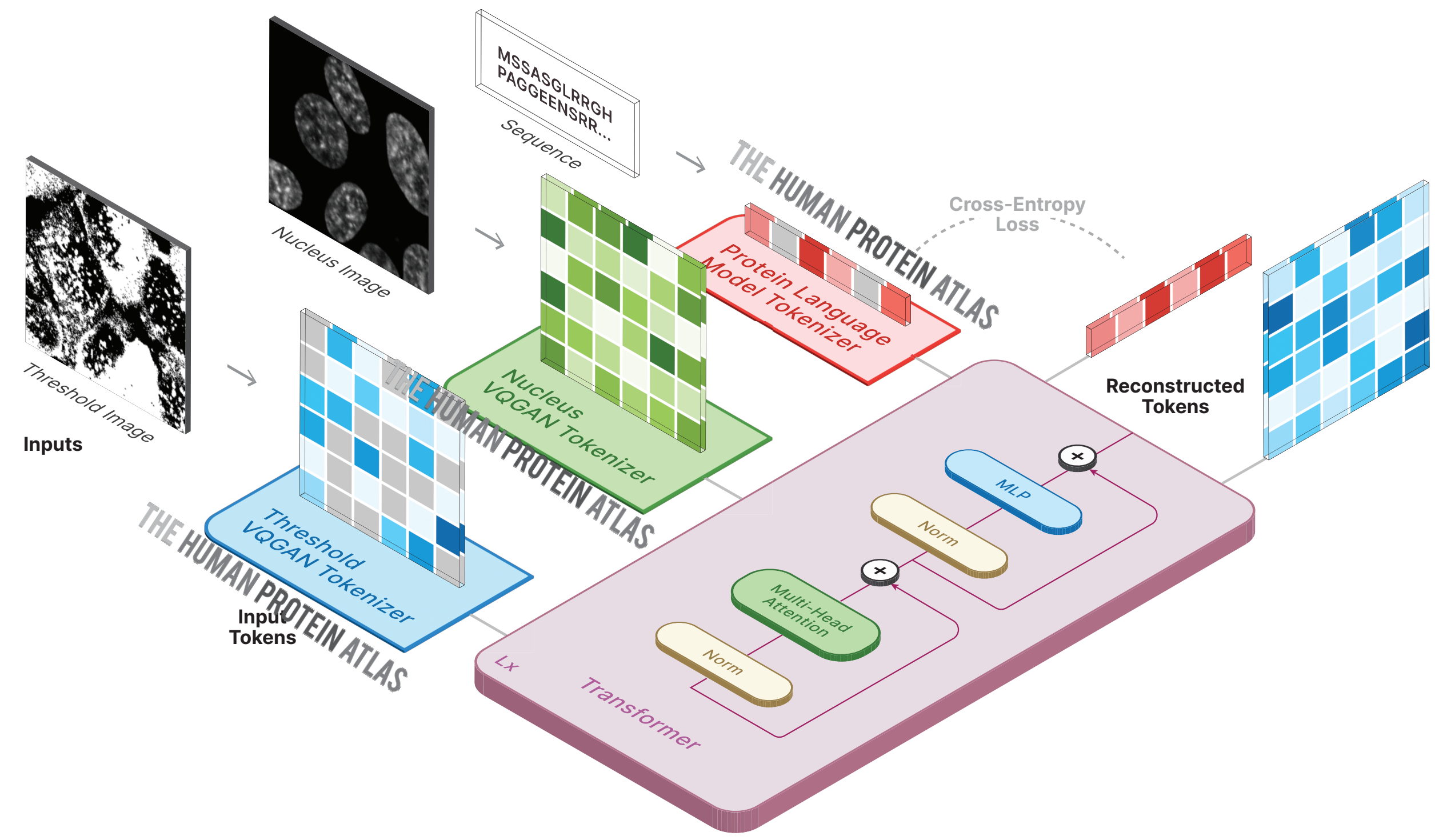


Finetuned Model Comparison Glycerol-3-phosphate acyltransferase 4

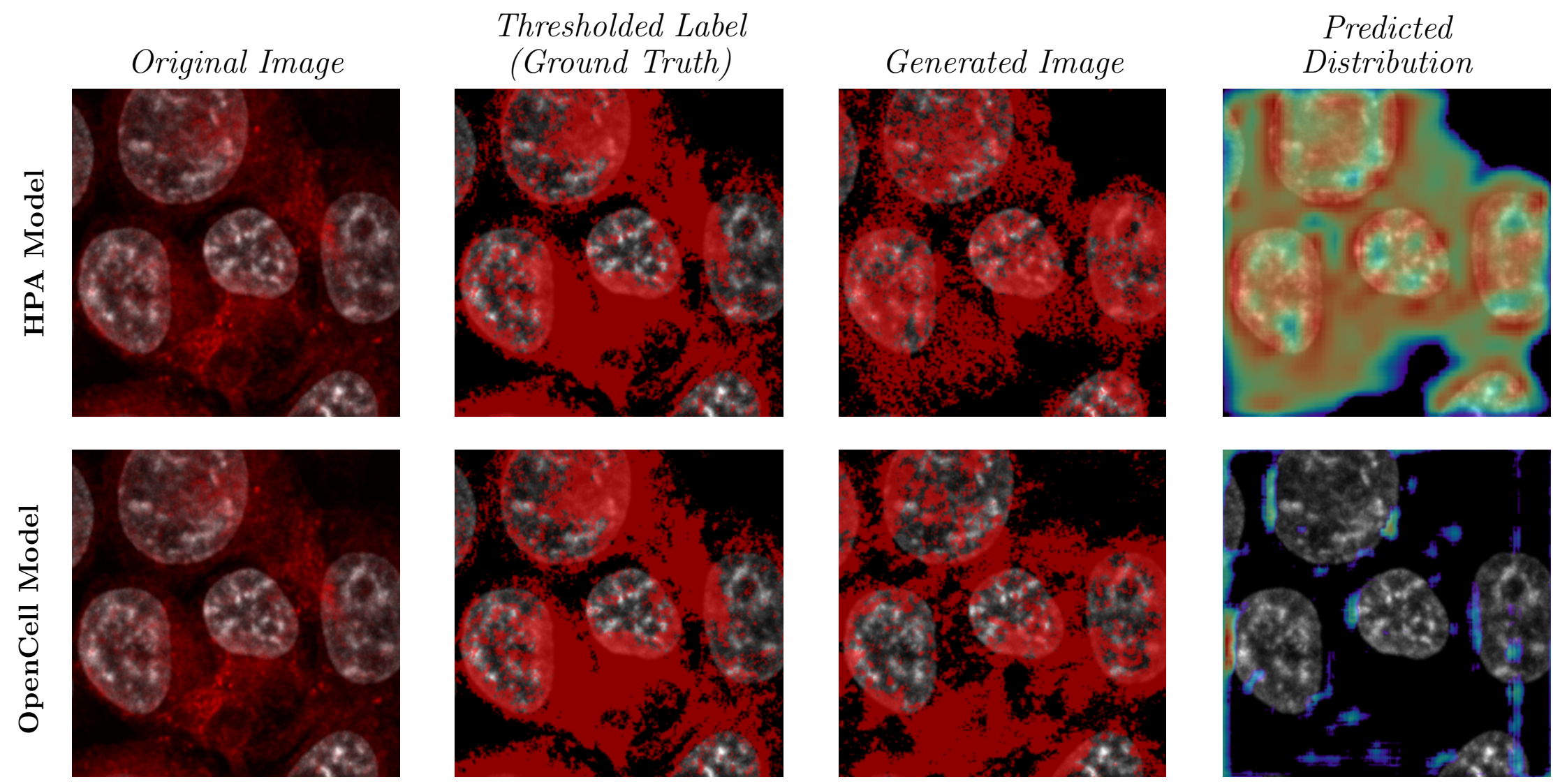
HPA Model



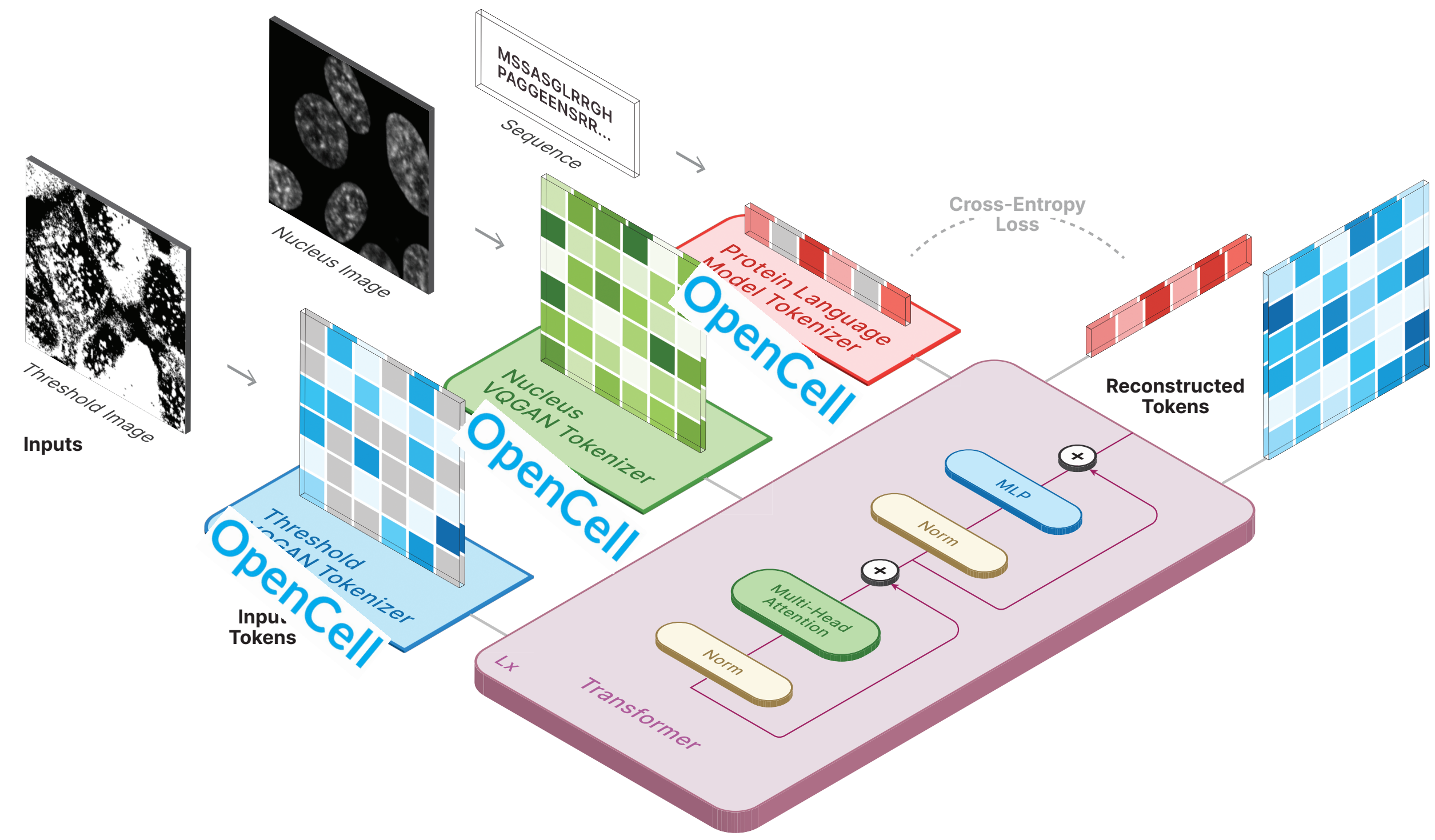
OpenCell Finetuning



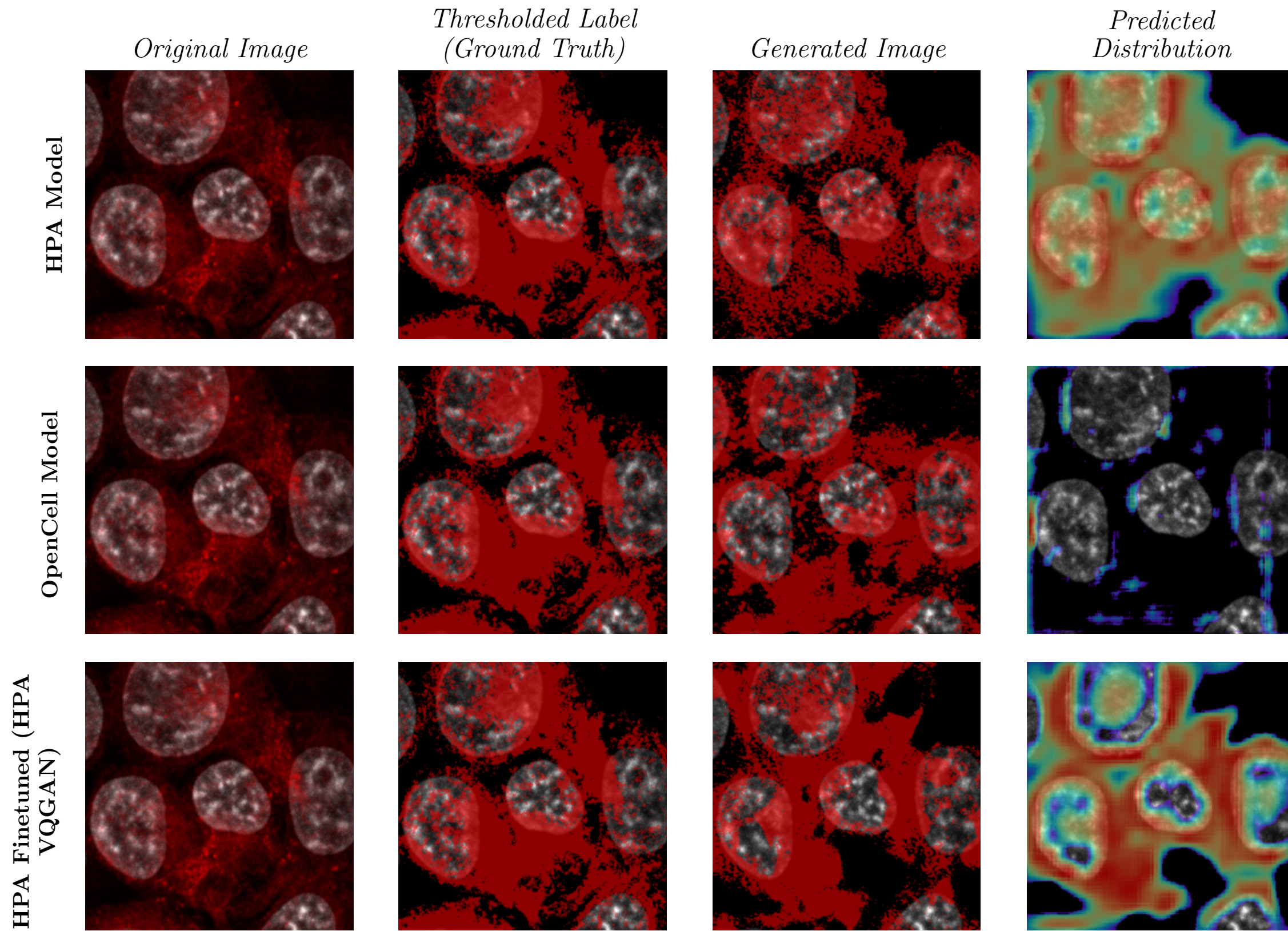
Finetuned Model Comparison Glycerol-3-phosphate acyltransferase 4



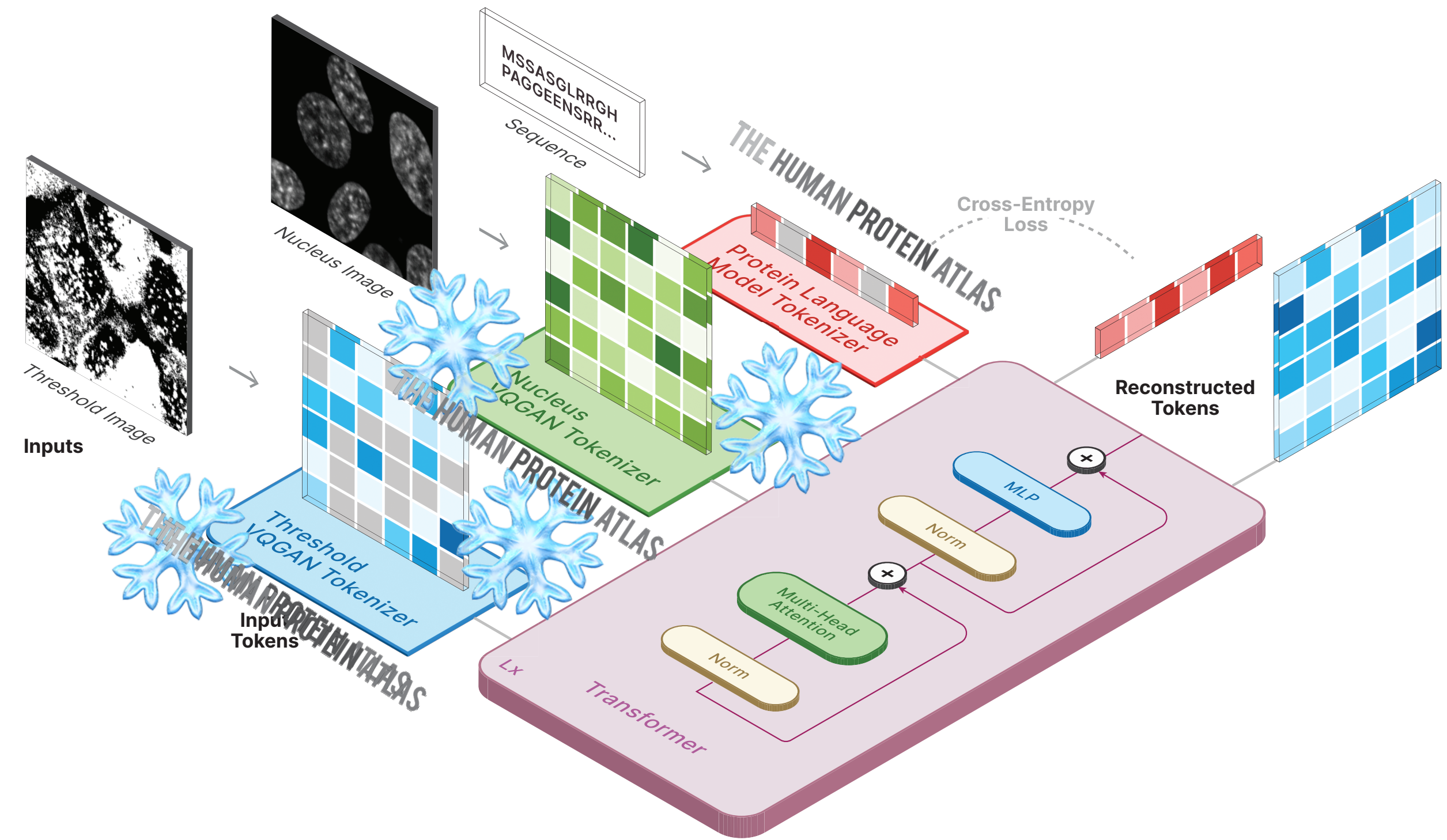
OpenCell Finetuning



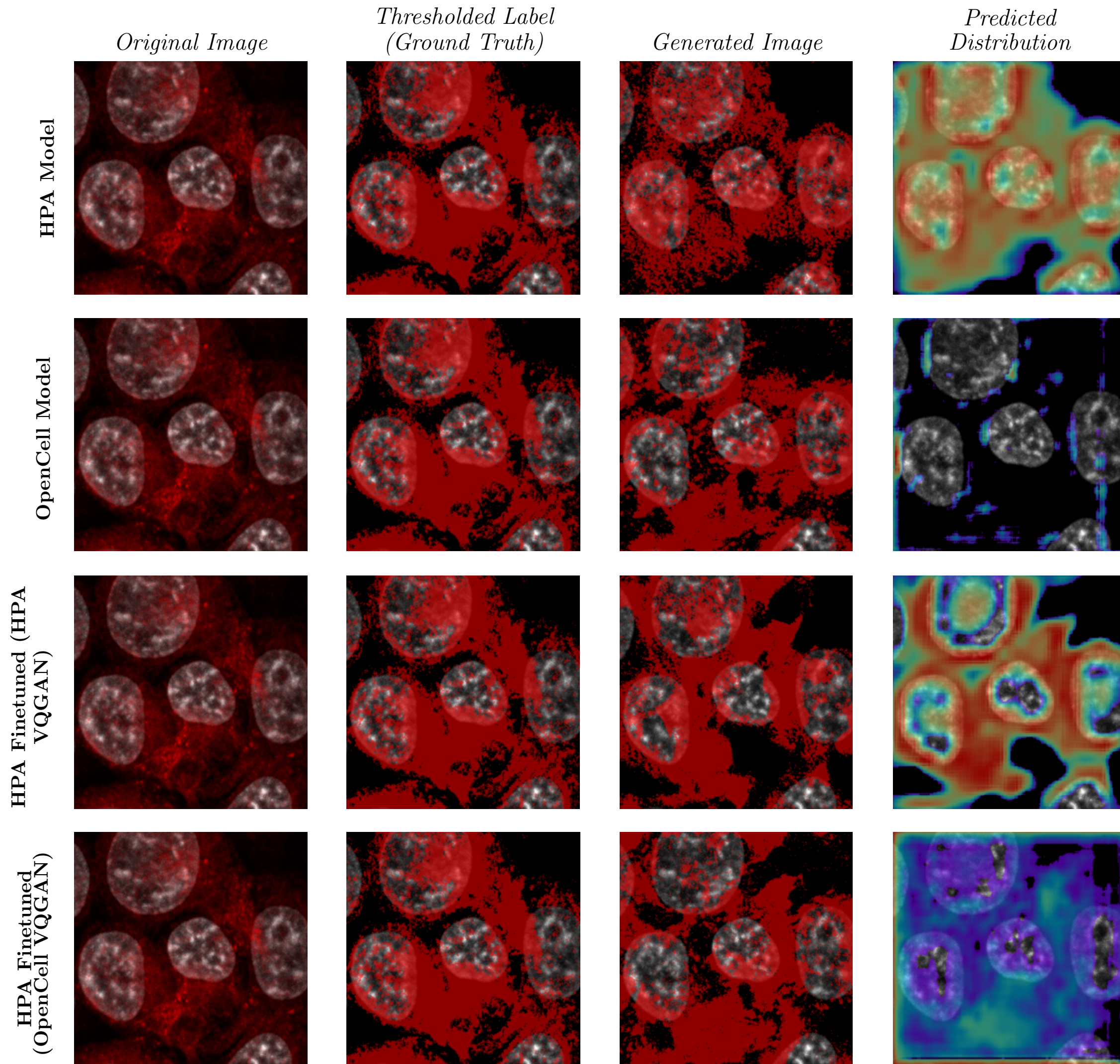
Finetuned Model Comparison
Glycerol-3-phosphate acyltransferase 4



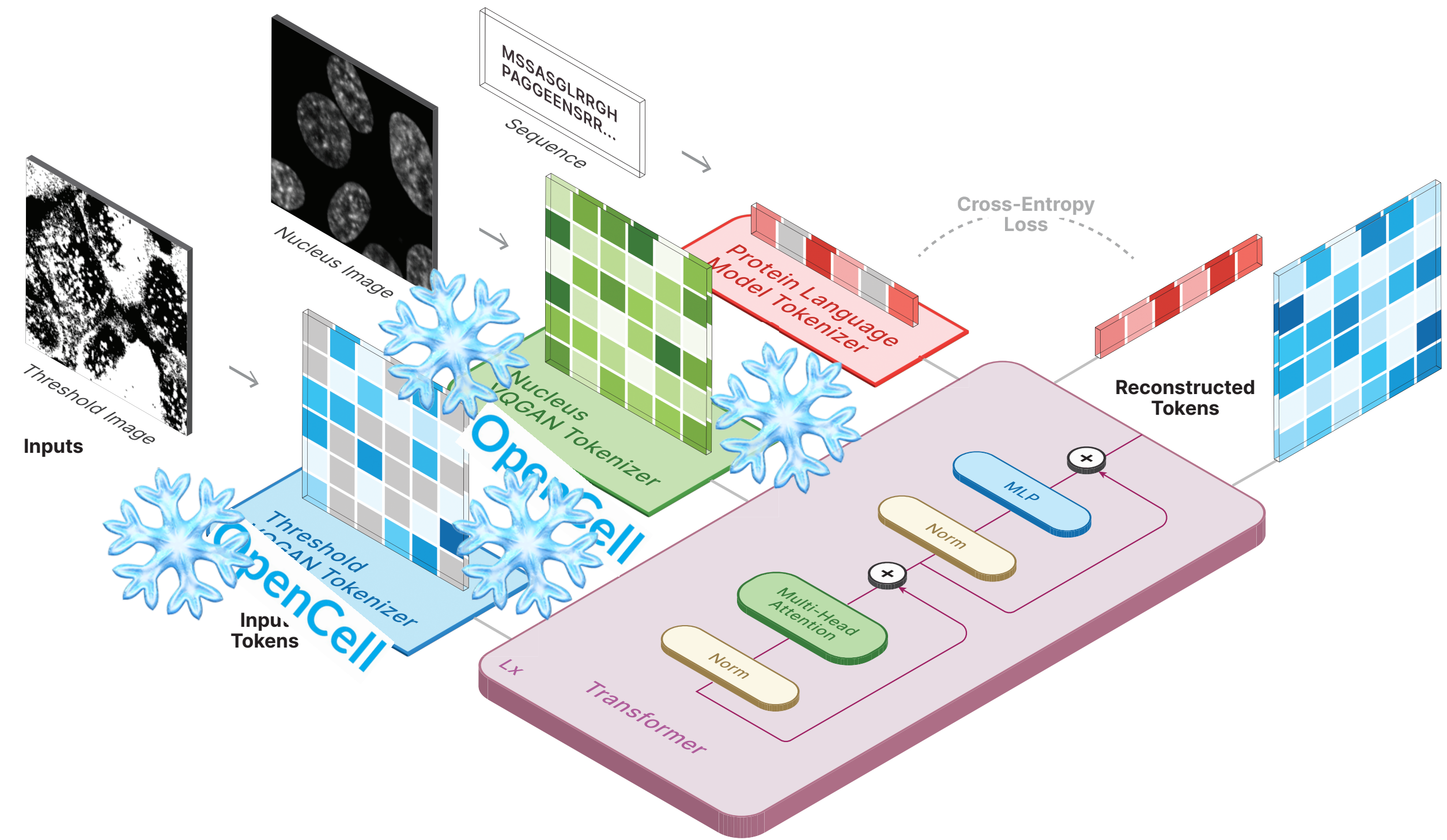
OpenCell Finetuning



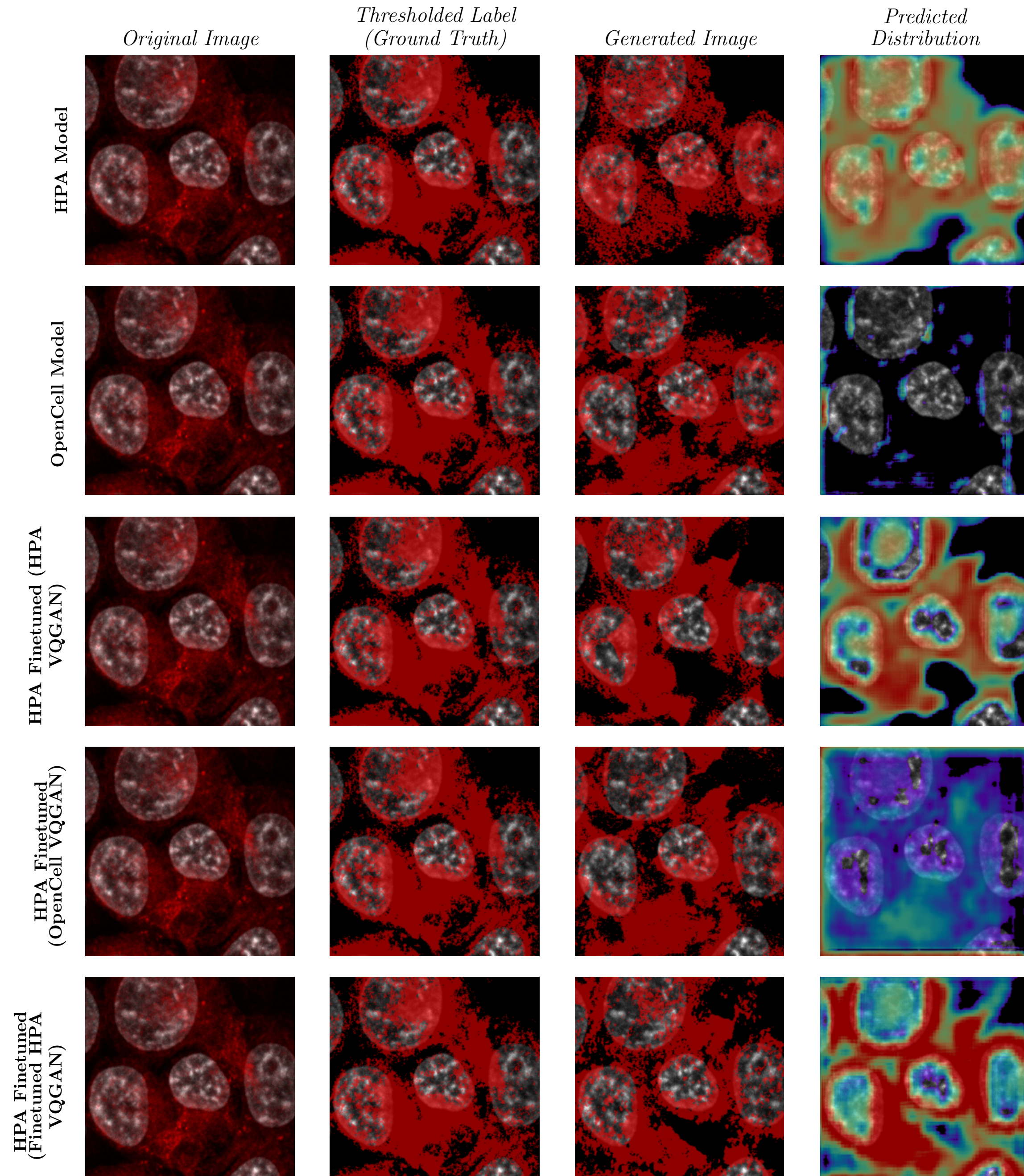
Finetuned Model Comparison
Glycerol-3-phosphate acyltransferase 4



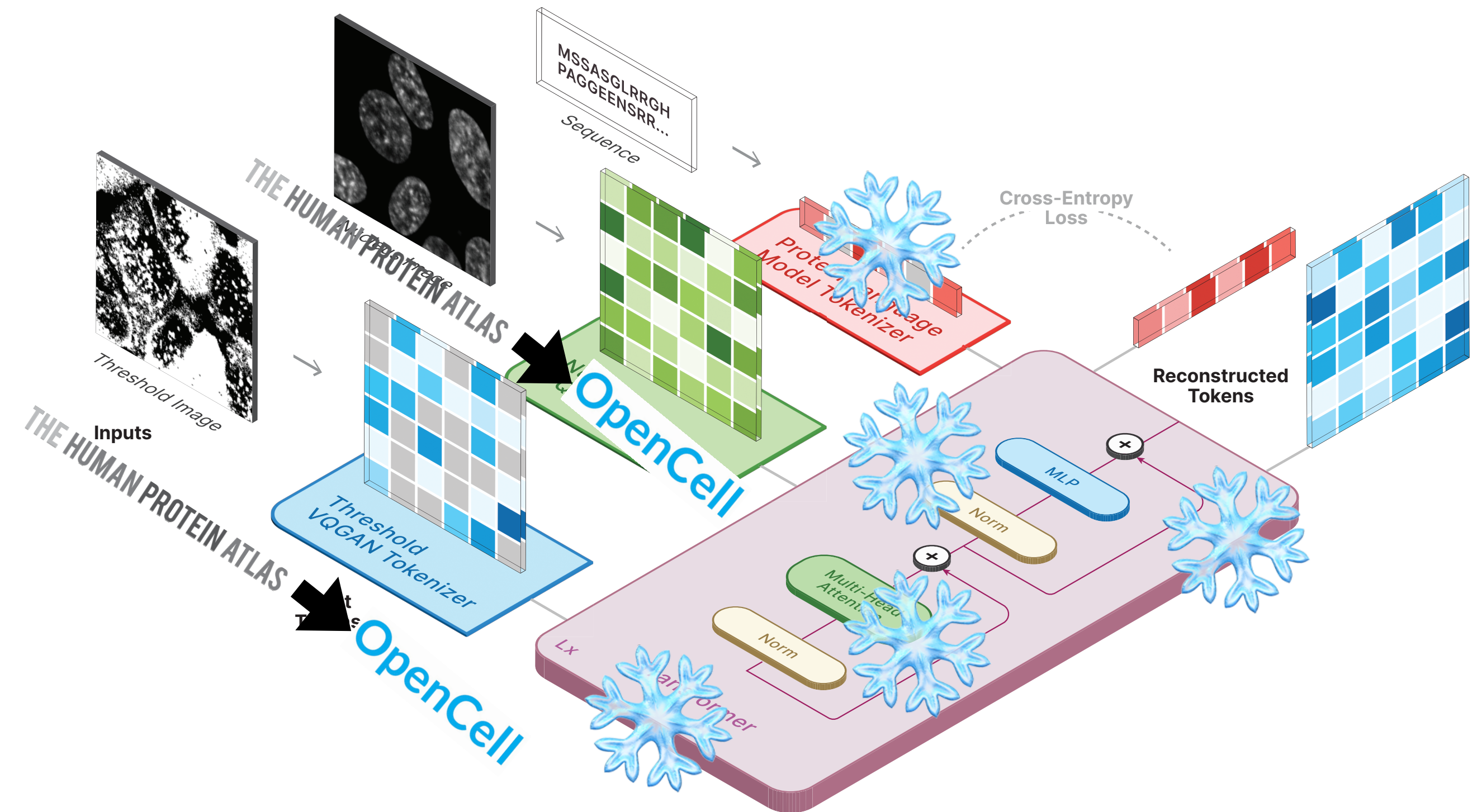
OpenCell Finetuning



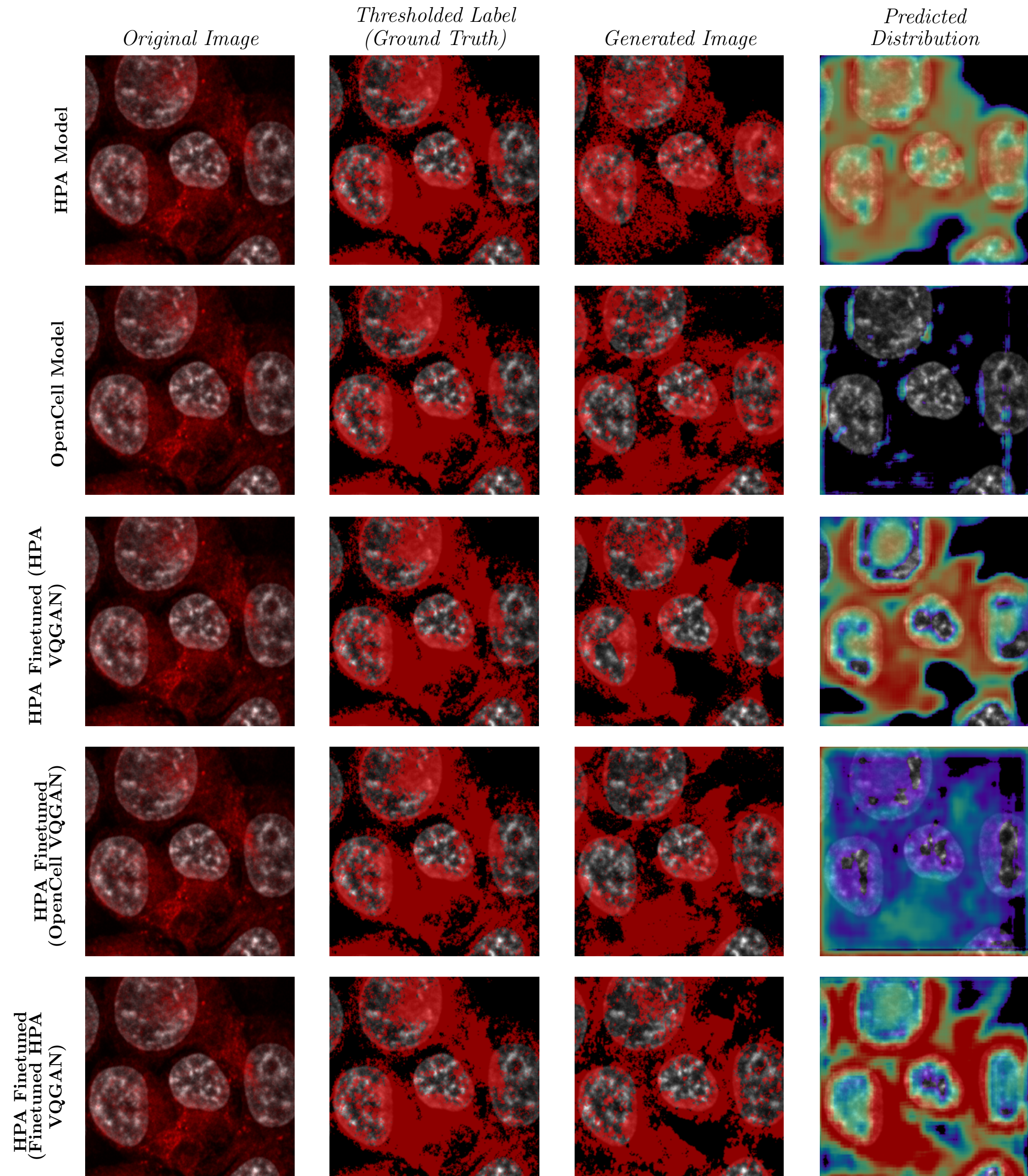
Finetuned Model Comparison
Glycerol-3-phosphate acyltransferase 4



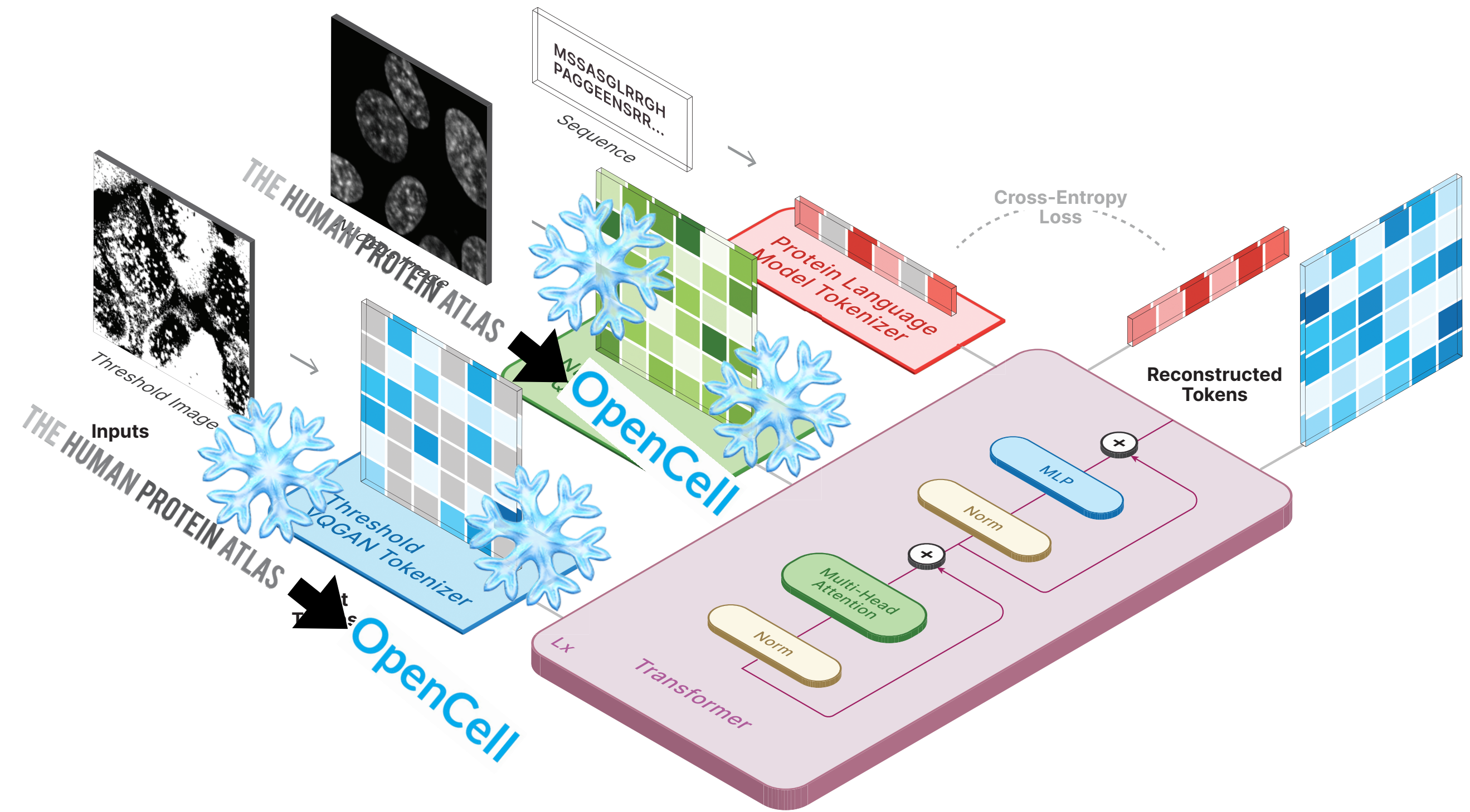
OpenCell Finetuning



Finetuned Model Comparison
Glycerol-3-phosphate acyltransferase 4



OpenCell Finetuning



CELL-E vs CELL-E 2

On OpenCell Validation Set

	Nucleus % MAE	Pixel Image MAE	PDF MAE	SSIM	FID	IS
CELL-E	.0347 ± .0294	.3671 ± .1117	.3653 ± .1008	.2060 ± .1846	10.5555	2.4762 ± .0866
CELL-E 2	.0170 ± .0160	.3449 ± .1305	.3487 ± .1340	.1881 ± .1541	19.2683	3.6083 ± .2013

Table S8: Speed Comparison

Model	Hidden Size	Autoregressive	Mean Generation Time (s)
CELL-E (Cached)	768	Yes	18.2740 \pm 0.0451
CELL-E (Non-Cached)	768	Yes	28.7694 \pm 0.3207
CELL-E 2	480	Yes	55.0057 \pm 0.2069
CELL-E 2	640	Yes	62.9650 \pm 0.1033
CELL-E 2	1280	Yes	74.3698 \pm 0.1788
CELL-E 2	2560	Yes	128.9960 \pm 0.3718
CELL-E 2	480	No	0.2784 \pm 0.0006
CELL-E 2	640	No	0.3067 \pm 0.0012
CELL-E 2	1280	No	0.3249 \pm 0.0011
CELL-E 2	2560	No	0.5487 \pm 0.0022

**66x Faster
with Better
Performance!**

CELL-E 1

Sequence

```
MSGTHLHNDLSQIEAIFRLNDSHKHKDKH  
KDREHRHKEHKKEKDREKSKHSNSEHKD  
SEKKEKKEKTKHKDGSSEKHKDKHKDR  
DKEKRKEEKVVRASGDAKIKKEKENGFS  
PPQIKDEPEDDGYFVPPKEIKPLKRPR  
DEDDADYKPKKIKTETDKKEKRRKLEEE  
EDGKLLKPKNKDKDKKVPEDNKKKPK  
KEEEQKWKWEEERYPEGIKWKFLHKG  
PVFAAPPYEPLEENVKFYDYGKVKLS  
AEEVATFFAKMLDHEYTTKEIFRKNFFK  
DWRKEMINEEKNIITNLSKCDFTQMSQY  
FKAQTEARKQMSKEEKLKIKEENEKLLK  
EYGFCIMDNHKERIANFKIEPPGLFRGR  
GNHPKMGMLKRRIMPEDIINCSKDAKV  
PSPPPGHKWKVVRHDNKVTWLVSWTENI  
QGSIKYIMLNPSSRIKGEKDWQKYETAR  
RLKKCVDKIRNQYREDWKSKEMKVRQRA  
VALYFIDKLLALRAGNEKEEGETADTVGC  
CSLRVEHINLHPELDGQEYVVEFDLFGK  
DSIRYYNKVPVEKRVFKNLQLFMENKQP  
EDDLFDRLNTGILNKHLQDLMEGLTAKV  
FRTYNASITLQQQLKELTAPDENIPAKI  
LSYNRANRAVAAILCNHQRAPPKTFEKSM  
MNLQTKIDAKKEQLADARRDLKSAKADA  
KVMKDAKTKKVVESKKAQVQRLEEQLMK  
LEVQATDREENKQIALGTSKLNYPRI  
TVAWCCKWGVPIEKIYNKTQREKFAWA  
DMADEDEYE
```

Nucleus Image



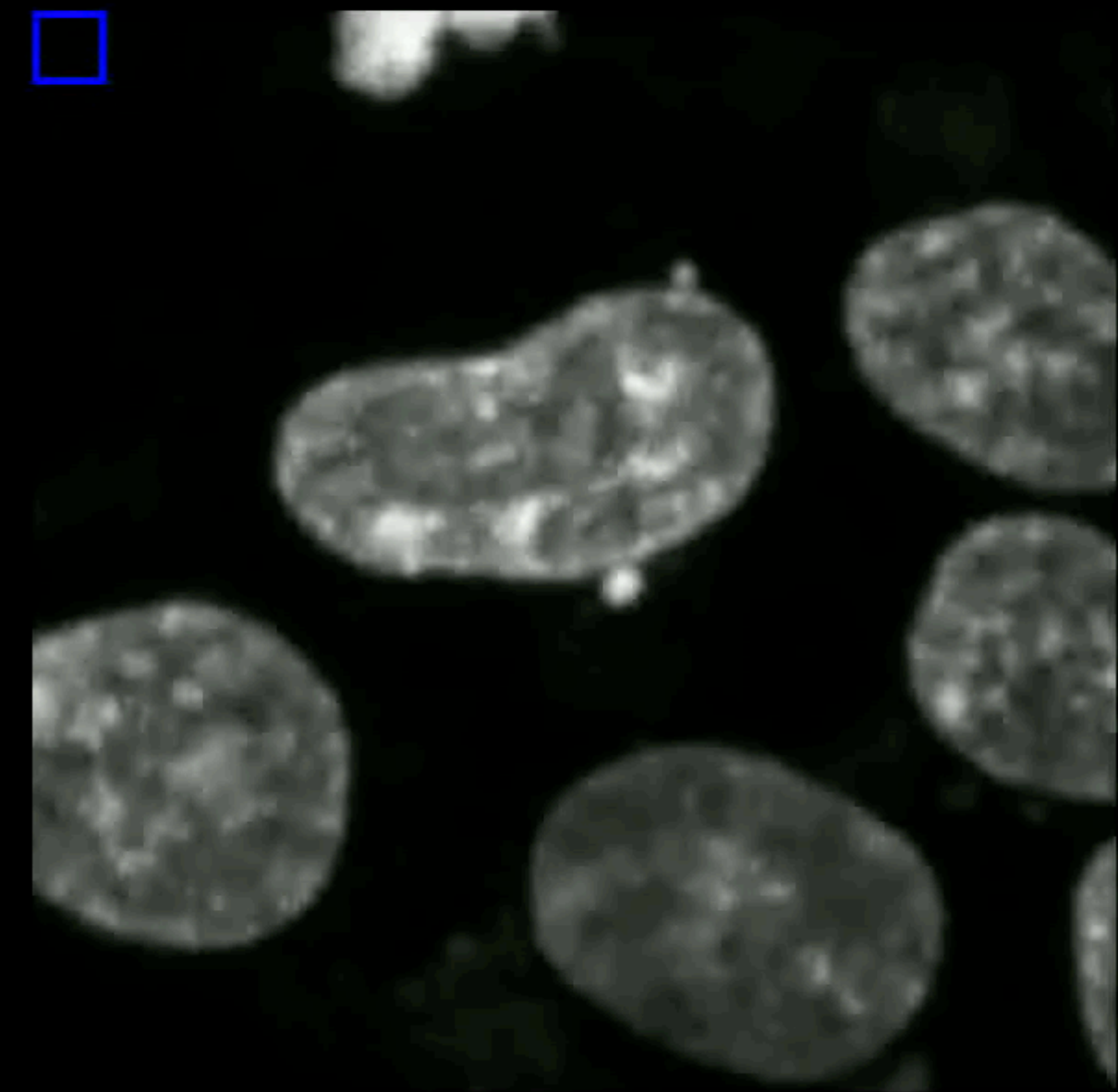
DNA Topoisomerase I
Relative Attention Weights



Predicted Threshold Image



Predicted PDF

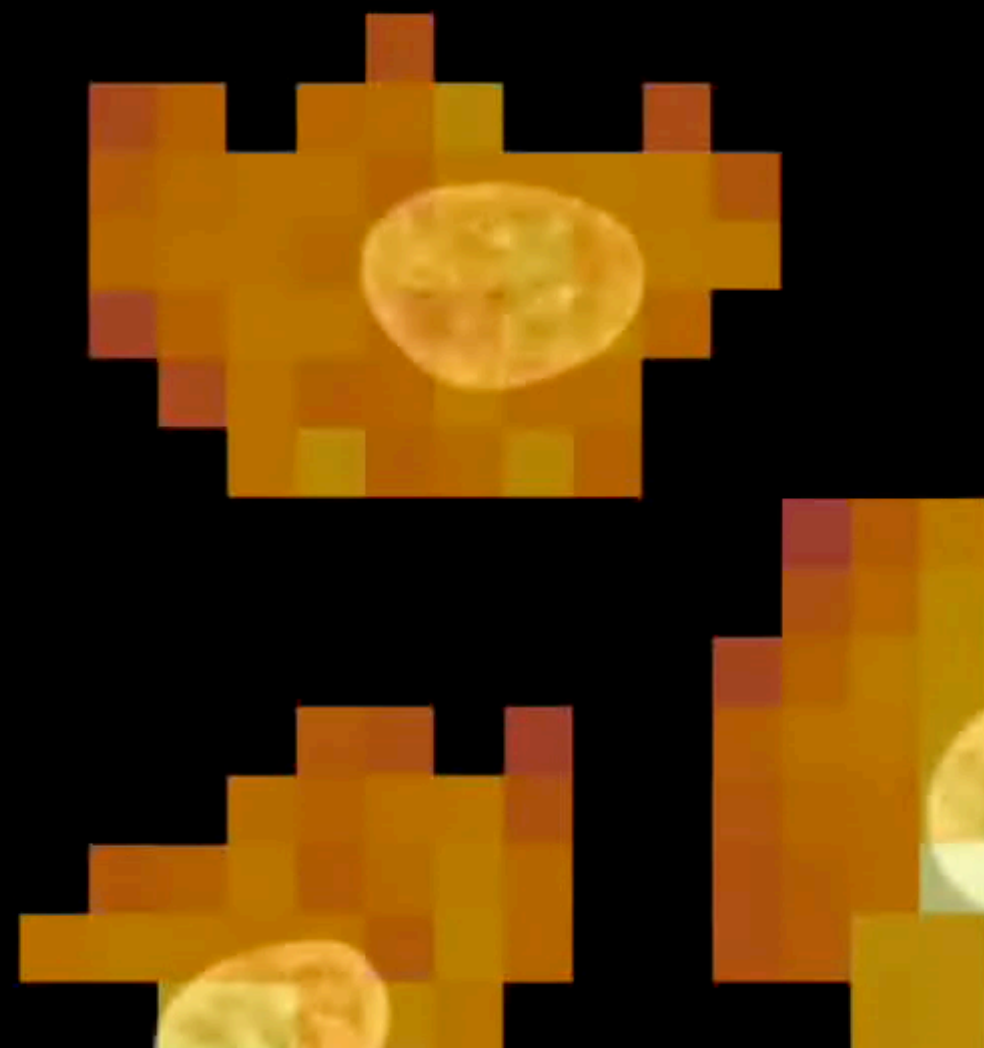


CELL-E 2

Sequence

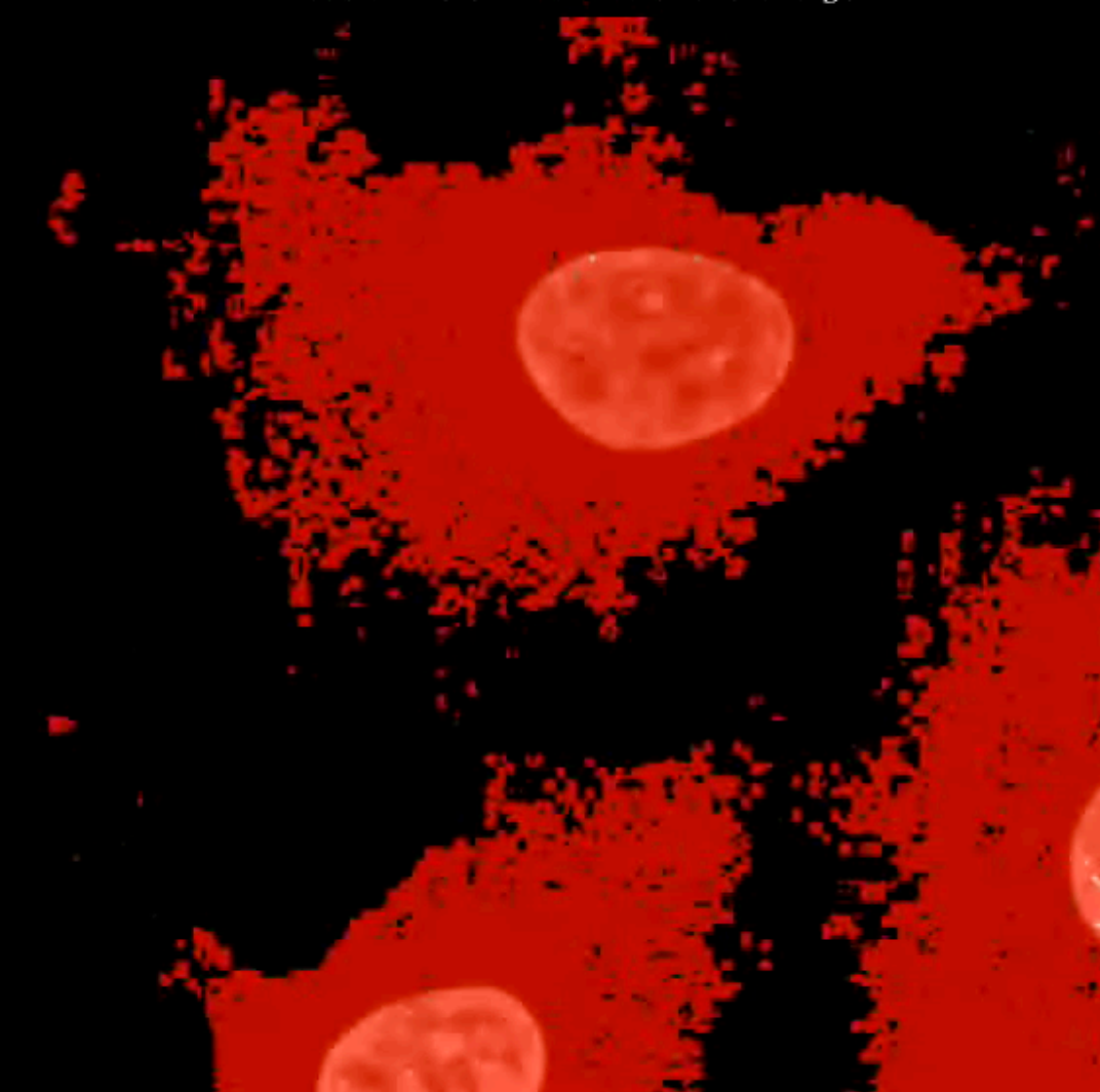
```
MSSPISKSRSLAFLQQLRSPRQ  
PPRLVTSSTAYTSPQPREVPVCP  
TAGGETQNAALPPTSWPLLGS  
LLQILWKGGLKKQHDTLVEYHKK  
YGKIFRMKLGSESVHLGSPCLL  
EALYRTESAYPQRLEIKPWKAYR  
DYRKEGYGLLILEGEDWQRVRS  
FQKKLMKPGEMKLDNKINEVLA  
DFMGRIDELCDERGHVEDLYSEL  
NKWSFESICLVLYEKRFGLLQKN  
AGDEAVNFIMAIKIMMSTFG  
VTPVELLHKS LNTKVWQDHTLAWD  
TIFKSVKACIDNRLLEKYSQQPSA  
DFLCDIYHQNRLSKKELYAAVTE  
LQLAAVETTANSLMWILYNLSRN  
PQVQQKLKELIQSVLPENQVPRA  
EDLRNMPYLKACLKESMRLTPSV  
PFTTRTLDKATVLGEYALPKGT  
LMLNTQVLGSSSEDNFEDSSQFRP  
ERWLQEKELKINPFAHLPPFVVGKR  
MCIGRRLAELQLHLALCWIVRKY  
DIQATDNEPVELMHSGLTLVPSRE  
LP I A F C Q R
```

Nucleus Image

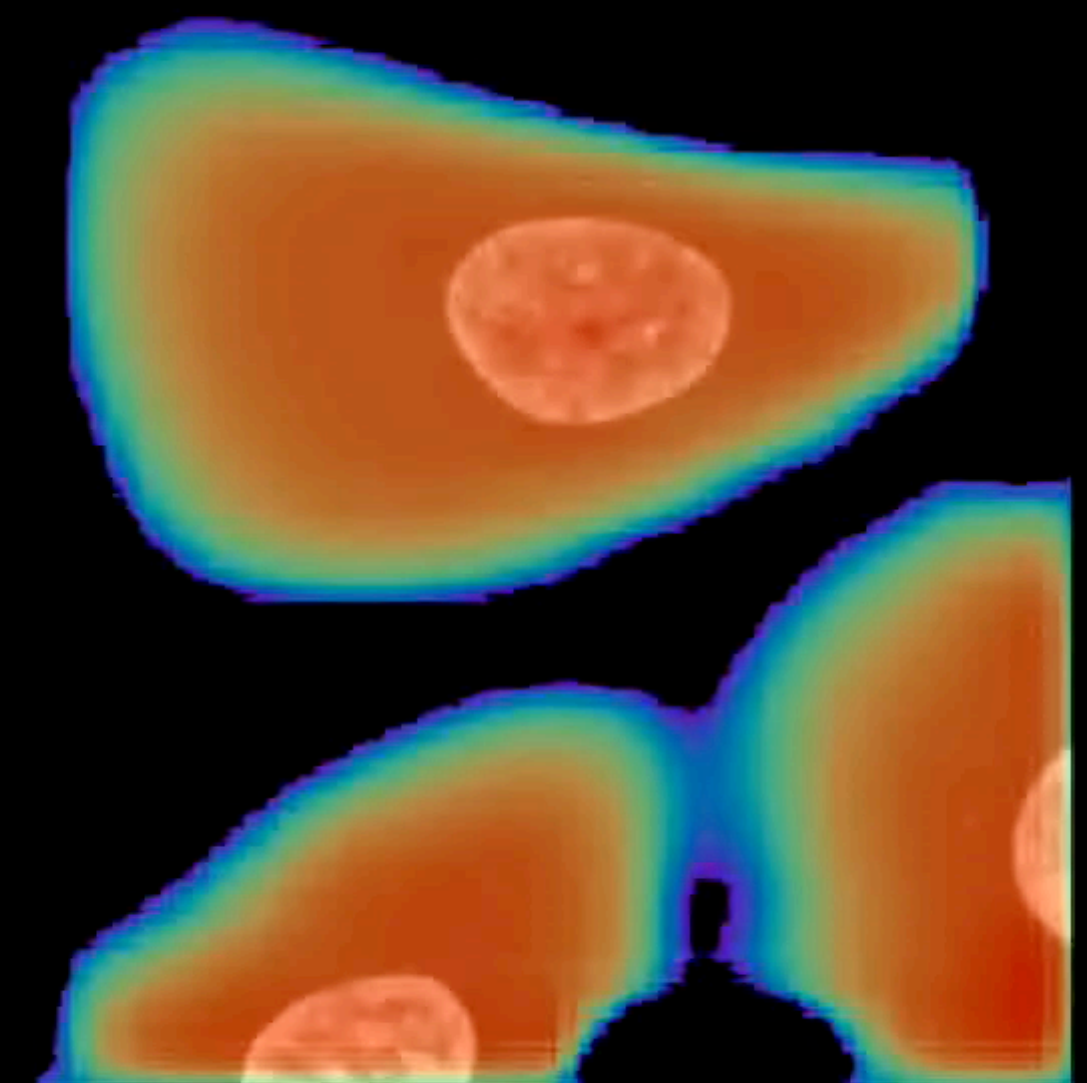


1,25-dihydroxyvitamin D(3) 24-hydroxylase, mitochondrial
Relative Attention Weights

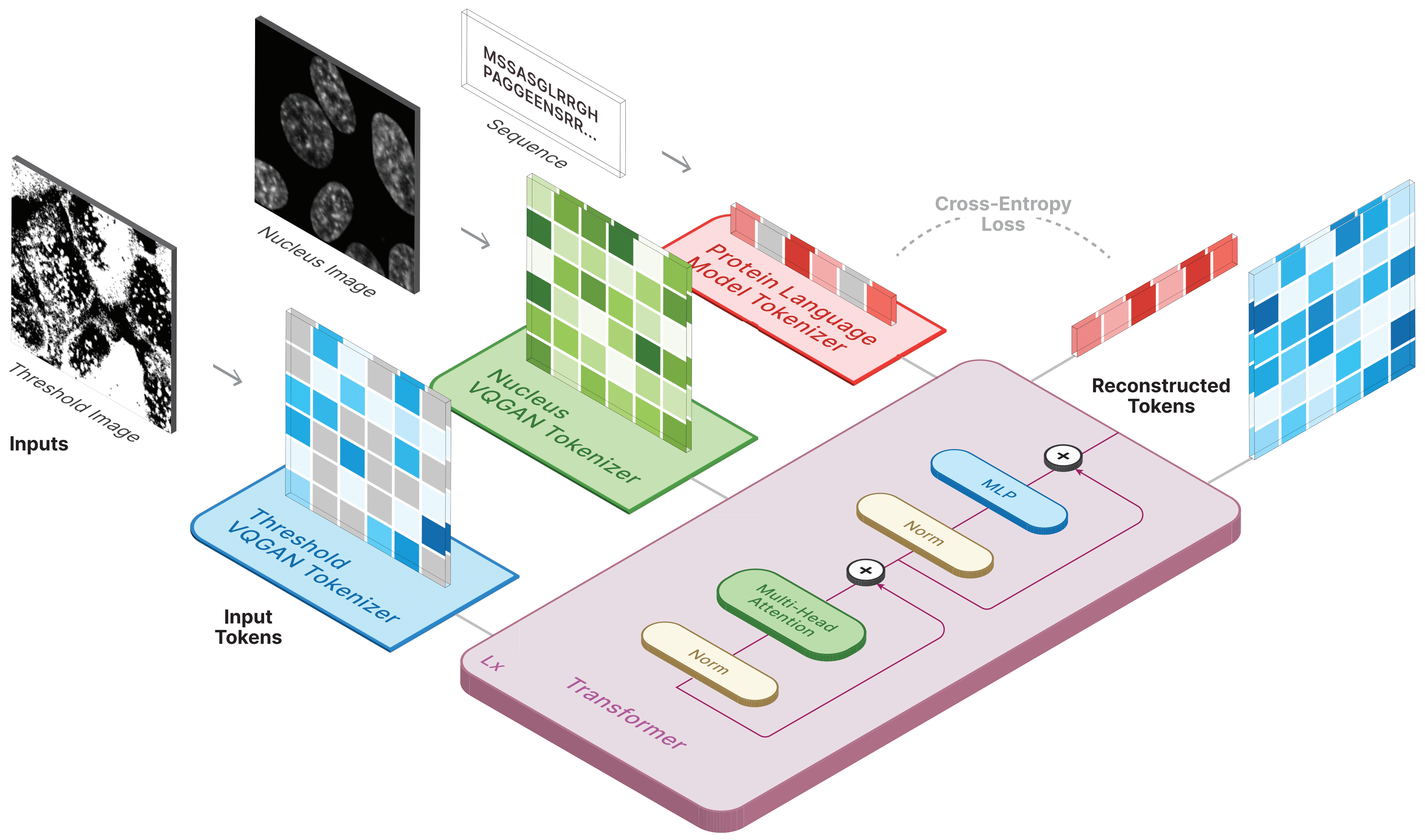
Predicted Threshold Image



Predicted PDF



De novo Protein Design

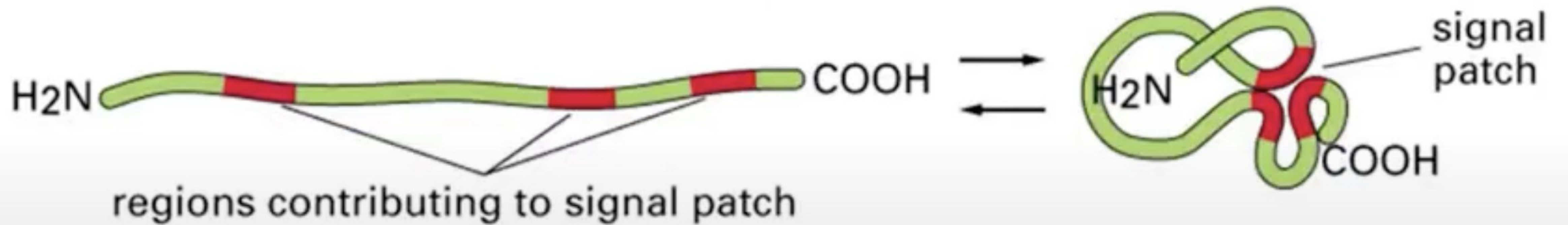


UNFOLDED PROTEIN

FOLDED PROTEIN



(A)



NLSdb sequence search results

Show **10** entries

Copy

CSV

Previous

1

Next

Query	Signal	SignalType	Start	End	ConfNuc	ConfFam	AnnotationType	Origin	Added	Modified
query	DKEKRRK	NLS	85	90	3	3	Potential	In Silico Mutagenesis	2017-10-11	2017-10-11
query	KEKKRK	NLS	159	164	6	6	Potential	In Silico Mutagenesis	2017-10-11	2017-10-11
query	KKRKLE	NLS	161	166	14	13	By Expert	P02545 , P48678 , P48679 , Q3ZD69	2017-10-11	2017-10-11
query	PLKRPR	NLS	135	140	4	3	Potential	In Silico Mutagenesis	2017-10-11	2017-10-11
All	All	All					All	All		

Showing 1 to 4 of 4 entries

Previous

1

Next

This is an updated version of NLSdb. The legacy NLSdb can be found [here](#). If you have questions and/or encounter any problems, you can contact us [here](#).
Copyright © 2017 Michael Bernhofer, Tatyana Goldberg and Burkhard Rost, [ROSTLAB](#) all rights reserved.

[CONTACT](#) | [HELP](#) | [IMPRINT](#)

Procedure

- Pre-select number of masked positions
- Generate 300 candidates per terminus with sequence model
- Feed predicted sequence to image model to cross-validate

MSKGEELFTGVVPILVELDGDVN
GHKFSVSGEGEGDATYGKLTLLKF
ICTTGKLPVPWPTLVTTFSYGVO
CFSRYPDHMKQHDFFKSAMPEGY
VQERTIFFKDDGNYKTRAEVKFE
GDTLVNRIELKGIDFKEDGNILG
HKLEYNYNSHNVYIMADKQKNGI
KVNFKIRHNIEDGSVQLADHYQQ
NTPIGDGPVLLPDNHYLSTQSAL
SKDPNEKRDHMLLEFVTAAGIT
HGMDELYK

Procedure

- Pre-select number of masked positions
- Generate 300 candidates per terminus with sequence model
- Feed predicted sequence to image model to cross-validate

```
<start>MSKGEELFTGVVPILV  
ELDGDVNGHKFSVSGEGEGDATY  
GKLTLLKFICTTGKLPVPWPTLVT  
TFSYGVQCFSRYPDHMKQHDFFK  
SAMPEGYVQERTIFFKDDGNYKT  
RAEVKFEGLTLVNRIELKGIDFK  
EDGNILGHKLEYNYNSHNVYIMA  
DKQKNGIKVNFKIRHNIEDGSVQ  
LADHYQQNTPIGDGPVLLPDNHY  
LSTQSALSKDPNEKRDHMLLEF  
VTAAGITHGMDELYK<end>
```

Procedure

- Pre-select number of masked positions
- Generate 300 candidates per terminus with sequence model
- Feed predicted sequence to image model to cross-validate

```
<start><mask>MSKGEELFTG  
VVPILVELDGDVNGHKFSVSGEG  
EGDATYGKLTCLKFICTTGKLPVP  
WPTLVTTFSYGVQCFSRYPDHMK  
QHDFFKSAMPEGYVQERTIFFKD  
DGNKTRAEVKFEGDTLVNRIEL  
KGIDFKEDGNILGHKLEYNNSH  
NVYIMADKQKNGIKVNFKIRHNI  
EDGSVQLADHYQQNTPIGDGPVL  
LPDNHYLSTQSALSKDPNEKRDH  
MVLLEFVTAAGITHGMDELYK<e  
nd>
```


Procedure

- Pre-select number of masked positions
- Generate 300 candidates per terminus with sequence model
- Feed predicted sequence to image model to cross-validate

```
<start><mask><mask><mask>  
k>MSKGEELFTGVVPILVELDGD  
VNGHKFSVSGEGEGDATYGKLT  
LKFICTTGKLPVPWPTLVTTFSY  
G  
VQCFSRYPDHMKQHDFFKSAMPE  
GYVQERTIFFKDDGNYKTRAEVK  
FEGDTLVNRIELKGIDFKEDGNI  
LGHKLEYNYNSHNVYIMADKQKN  
GIKVNFKIRHNIEDGSVQLADHY  
QQNTPIGDGPVLLPDNHYLSTQS  
ALSKDPNEKRDHMLLEFVTAAG  
ITHGMDELYK<end>
```

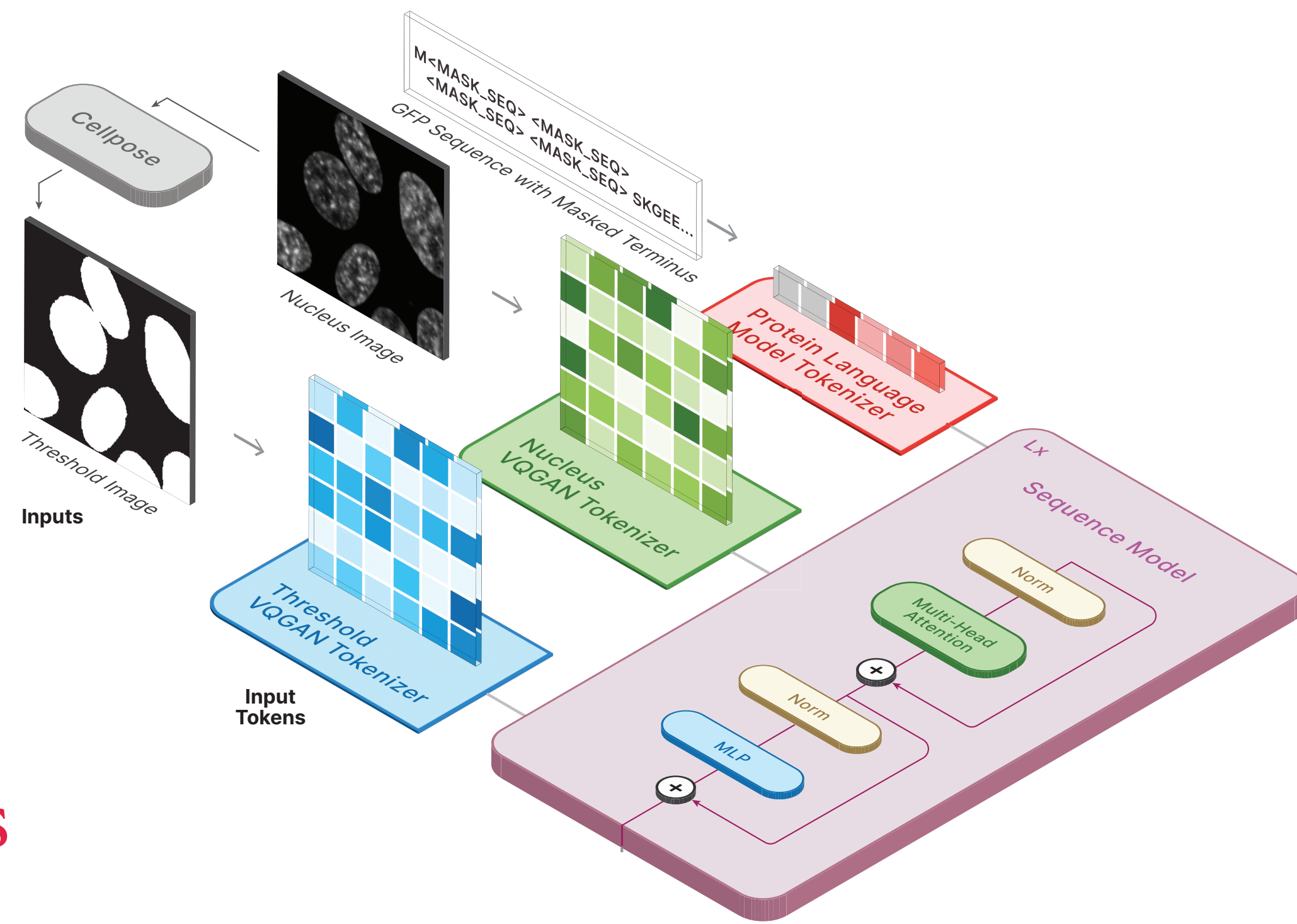
Procedure

- Pre-select number of masked positions
- Generate 300 candidates per terminus with sequence model
- Feed predicted sequence to image model to cross-validate

```
<start>MSKGEELFTGVVPILV  
ELDGDVNGHKFSVSGEGEGDATY  
GKLTLLKFICTTGKLPVPWPTLVT  
TFSYGVQCFSRYPDHMKQHDFFK  
SAMPEGYVQERTIFFKDDGNYKT  
RAEVKFEGLTLVNRIELKGIDFK  
EDGNILGHKLEYNYNSHNVYIMA  
DKQKNGIKVNFKIRHNIEDGSVQ  
LADHYQQNTPIGDGPVLLPDNHY  
LSTQSALSKDPNEKRDHMLLEF  
VTAAGITHGMDELYK<mask><m  
<ask><mask><end>
```


Procedure

- Pre-select number of masked positions
- Generate 300 candidates per terminus with sequence model
- Feed predicted sequence to image model to cross-validate



Procedure

- Pre-select number of masked positions
- Generate 300 candidates per terminus with sequence model
- Feed predicted sequence to image model to cross-validate

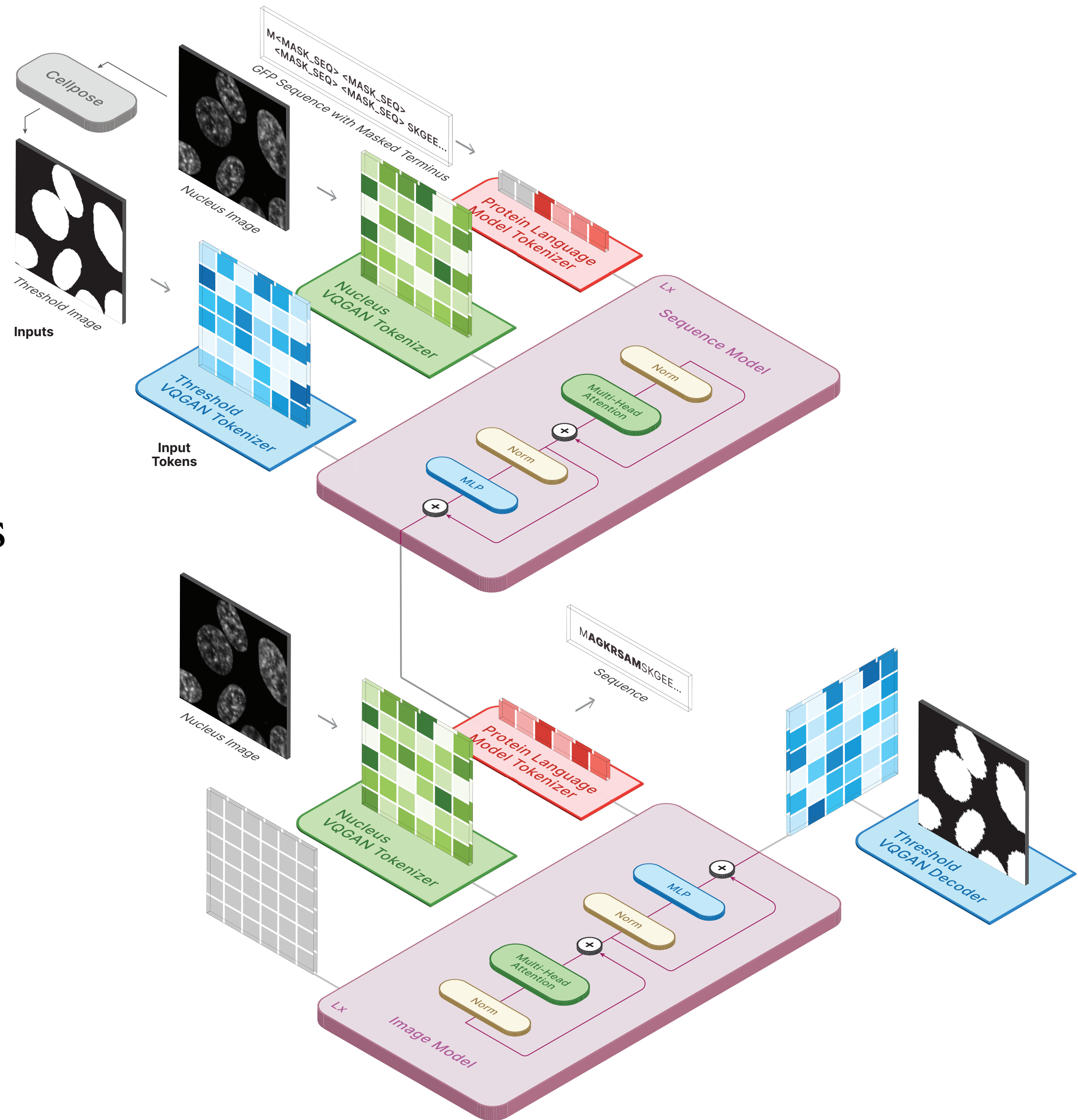
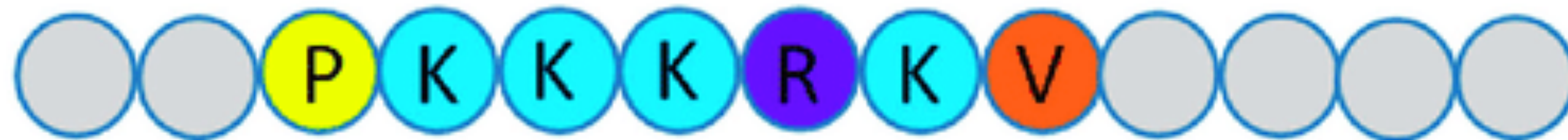


Table S10: NLS candidates sorted by nucleus proportion.

Terminus	Sequence	Terminus	Sequence
N	RKRRQR	C	SPTAFPSNVIETIRVKRRMEL
N	NKRPRKKEK	C	EFRAKYRQMGSRKKKKSGQWSA
C	RPKVI	N	KKHKLRSVPDLTELMRMIFLAP
C	VLKRAKGD	N	KLLRFAGKSGMMVLLAPHSGKM
C	RHKKKKIA	C	IFQADKDQKAHPPAKKAPSELMQ
N	HRRKKR	C	KGKVKSIMIPPKSRKSLAKVPLS
C	RSQKRK	N	AAGKSFKPRIKSRMTRDSSETMA
N	KCKKKN	C	TGNRIFGETPSWERERKRPGGGQQ
N	KGKRFSK	C	NKLQKHSKRQPHKLQAMKLKYPTWE
C	AKRLKGK	C	LVFPNRDASIKKPLQNPPQKRRCMIM
C	SKKAKKNKM	N	LPKRRRLSRRKKVELEPEYGWEEEVVV
C	EEKRPRF	N	TEAPARTAVKKS RAMKGYIARLASSPS
N	MKICIT	C	IEKSKGKEAPKSSPPLKQNQRSRKMVK
N	AVPAKRARIDG	C	FQVRASPKGKPATKNKLRLKIRRHRV
C	ESHHLPRAKKR	C	LQEGTRTRSQKAQEPKFKKVSGDIPNK
N	GKERSYPPISKR	N	SDPNTAQYPWMPPQATKRAAMAAREAE
C	KLKKRNRQPEDKK	C	HYKKEKRKRSASPILAEPPVKCARTLR
C	GGKFATGKKKKPKM	C	LDKRKRIKPPKEEQKELMRKMWGPGSSL
N	PSKLLRQ	N	GSKKSRTATDSLESRMAMEDVAMGEESE
C	QRRKGQKFQT	C	EGSGLVPGNSRKRPEPKPKKRKKVRRK
C	KTCPPKRPVVEW	C	RKKRQAIQAVTMGRIKKKSYEKQWSKFED
C	DKEKKRKNDHEK	C	ASTVPAYSRSKAGKVEPKPKQKKTQRNAP
N	FRFSC	C	SKOQAEINLKAAPLETTDISLSKKEKKDM

NLS Properties

- Short (5-15 Amino Acids)
- Highly Basic (R and K)
- Clustered



Monopartite NLSs



Bipartite NLSs

Max ID %	# Sequences	Mean Sequence Length	Mean % R or K
0% - 33%	109	25.6606 \pm 3.0099	20.6379 \pm 8.6101
33% - 66%	133	17.1955 \pm 5.0804	32.0076 \pm 12.8334
66% - 100%	13	6.9231 \pm 1.2558	57.5794 \pm 17.9351

Highly Dissimilar

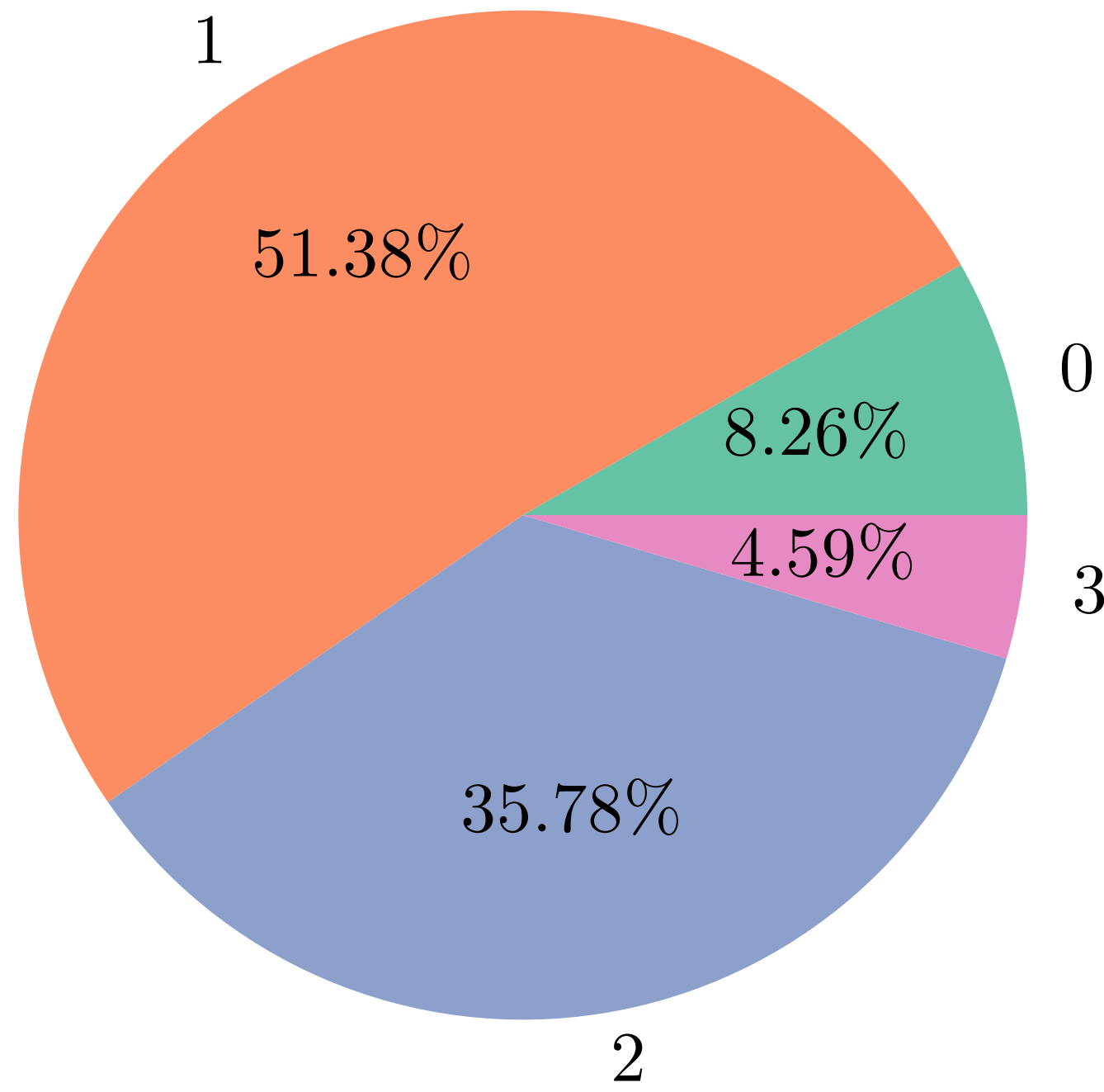


Max ID %	# Sequences	Mean Sequence Length	Mean % R or K
0% - 33%	109	25.6606 ± 3.0099	20.6379 ± 8.6101
33% - 66%	133	17.1955 ± 5.0804	32.0076 ± 12.8334
66% - 100%	13	6.9231 ± 1.2558	57.5794 ± 17.9351

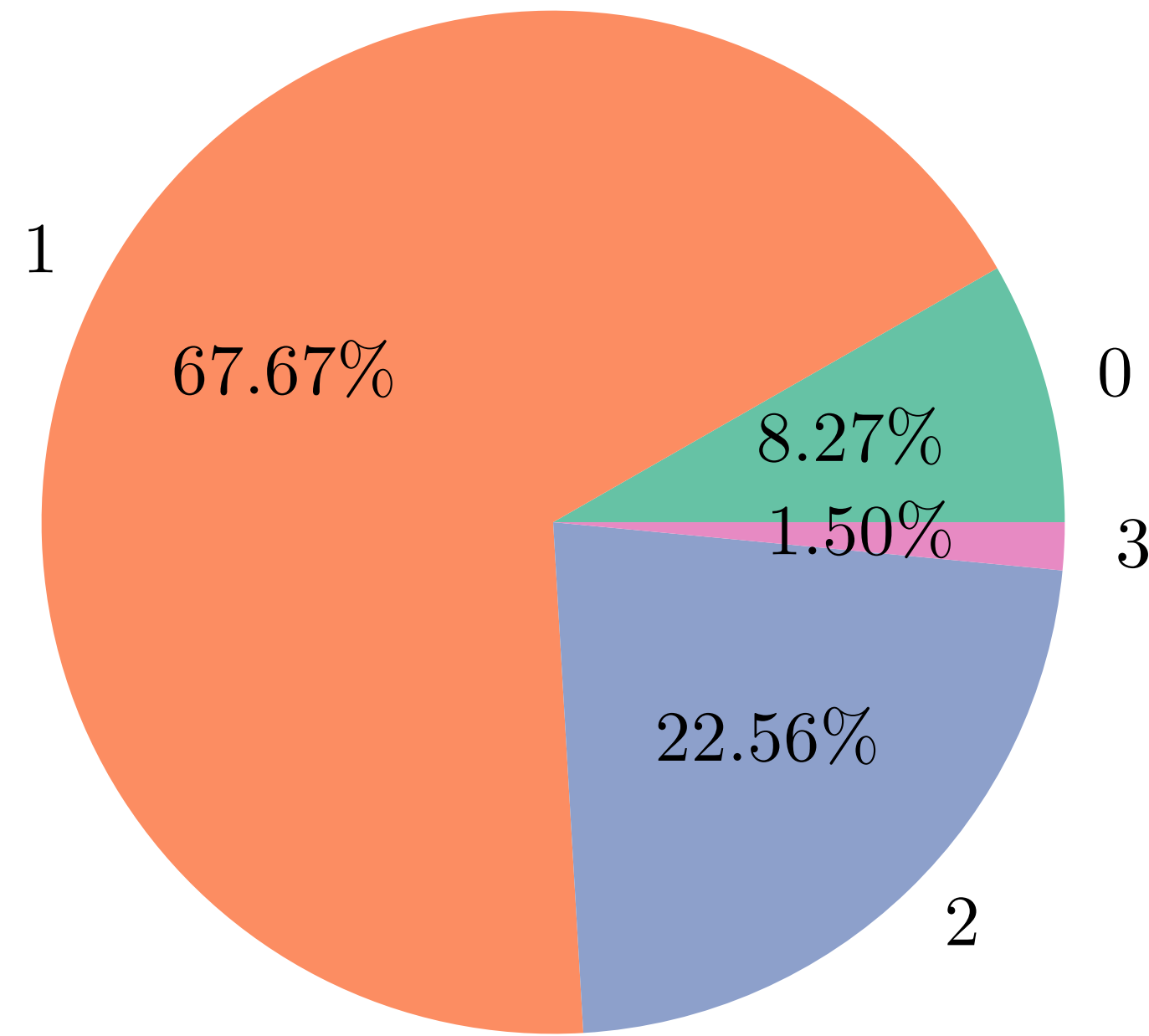
of Basic Amino Acid Stretches per Sequence

Max ID%

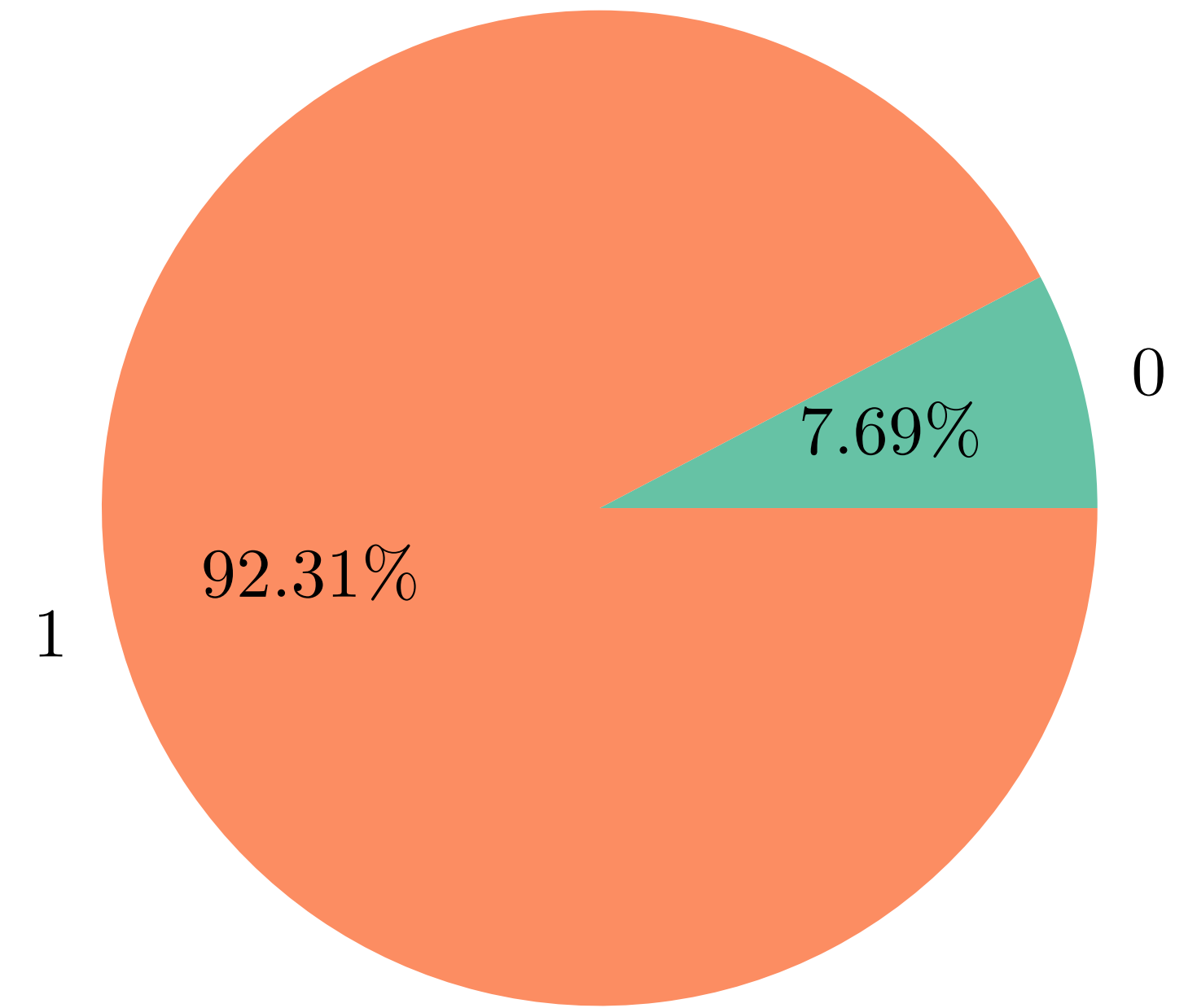
0% – 33%



33% – 66%



66% – 100%



DeepLoc 2.0 Prediction:

- 89% Nuclear Localizing Proteins**
- 91% contain NLS**

Future Directions

- Leverage more data
- Incorporate structural information
- Validate *de novo* candidates

Thanks!