



# MSc Thesis Committee Meeting # 1

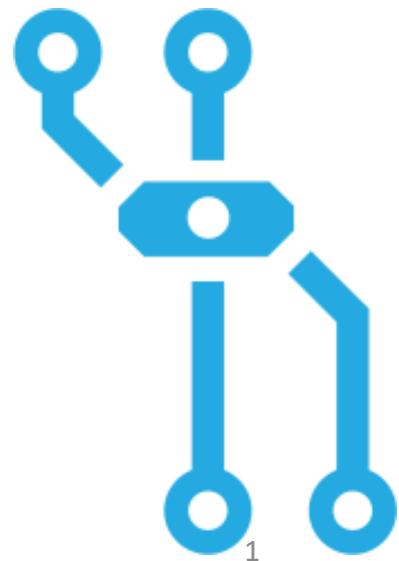
Monday, June 15, 2020



**Diana Lin**

MSc Student, Bioinformatics Graduate Program, UBC  
Dr. Inanc Birol Lab, Genome Sciences Centre (GSC), BC Cancer

Supervisor: Dr. Inanc Birol



# Committee Members



THE UNIVERSITY  
OF BRITISH COLUMBIA



**Dr. Inanc Birol, PhD**

Medical Genetics, UBC  
Birol Lab, Genome Sciences Centre, BC Cancer



**Dr. Phil Hieter, PhD**

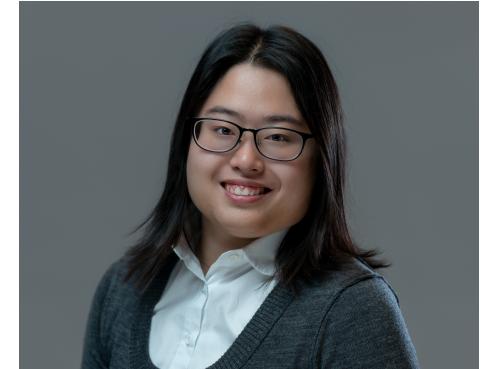
Medical Genetics, UBC  
Hieter Lab, Michael Smith Laboratories, UBC



**Dr. Michael Murphy, PhD**

Microbiology and Immunology, UBC  
Murphy Lab, Life Sciences Institute, UBC

# Agenda



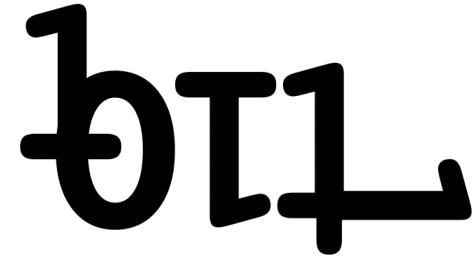
1. Introduction (3 min.)
2. Background (6 min.)
3. Progress Update (13 min.)
4. Open Discussion and Feedback (30 min.)

# Education (2014 - 2018)



- Bachelor of Science, Major Physiology, McGill University
  - Elective courses to supplement Life Science degree:
    - COMP 202: Foundations of Programming (Java)
    - COMP 206: Introduction to Software Systems (Unix CLI, C)
    - COMP 250: Introduction to Computer Science (pseudocode)
    - COMP 322: Introduction to C++ (C++)
    - COMP 364: Computer Tools for Life Science (Python)

# Internship (2018 - 2019)



- Student Intern, Birol Lab, Genome Sciences Centre
  - Data mining for antimicrobial peptide sequences in annotated nuclear spruce genomes using homology search
  - Differential expression analysis of AMPs in spruce transcriptomes
  - Assembly and annotation of organellar genomes

Species	Genome	Assembly	Annotation	NCBI Accession
White Spruce	chloroplast		✓	MK174379.1
Engelmann Spruce	chloroplast	✓	✓	MK241981.1
White Pine Weevil	mitochondrion	✓	✓	preparing submission



## Complete Chloroplast Genome Sequence of a White Spruce (*Picea glauca*, Genotype WS77111) from Eastern Canada

Diana Lin,<sup>a</sup> Lauren Coombe,<sup>a</sup> Shaun D. Jackman,<sup>a</sup> Kristina K. Gagalova,<sup>a</sup> René L. Warren,<sup>a</sup> S. Austin Hammond,<sup>a\*</sup> Heather Kirk,<sup>a</sup> Pawan Pandoh,<sup>a</sup> Yongjun Zhao,<sup>a</sup> Richard A. Moore,<sup>a</sup> Andrew J. Mungall,<sup>a</sup> Carol Ritland,<sup>b,f</sup> Barry Jaquish,<sup>c</sup> Nathalie Isabel,<sup>d</sup> Jean Bousquet,<sup>e</sup> Steven J. M. Jones,<sup>a</sup> Joerg Bohlmann,<sup>b,f</sup> Inanc Birol<sup>a</sup>

## Complete Chloroplast Genome Sequence of an Engelmann Spruce (*Picea engelmannii*, Genotype Se404-851) from Western Canada

Diana Lin,<sup>a</sup> Lauren Coombe,<sup>a</sup> Shaun D. Jackman,<sup>a</sup> Kristina K. Gagalova,<sup>a</sup> René L. Warren,<sup>a</sup> S. Austin Hammond,<sup>a\*</sup> Helen McDonald,<sup>a</sup> Heather Kirk,<sup>a</sup> Pawan Pandoh,<sup>a</sup> Yongjun Zhao,<sup>a</sup> Richard A. Moore,<sup>a</sup> Andrew J. Mungall,<sup>a</sup> Carol Ritland,<sup>b,e</sup> Trevor Doerksen,<sup>c</sup> Barry Jaquish,<sup>c</sup> Jean Bousquet,<sup>d</sup> Steven J. M. Jones,<sup>a</sup> Joerg Bohlmann,<sup>b,e</sup> Inanc Birol<sup>a</sup>

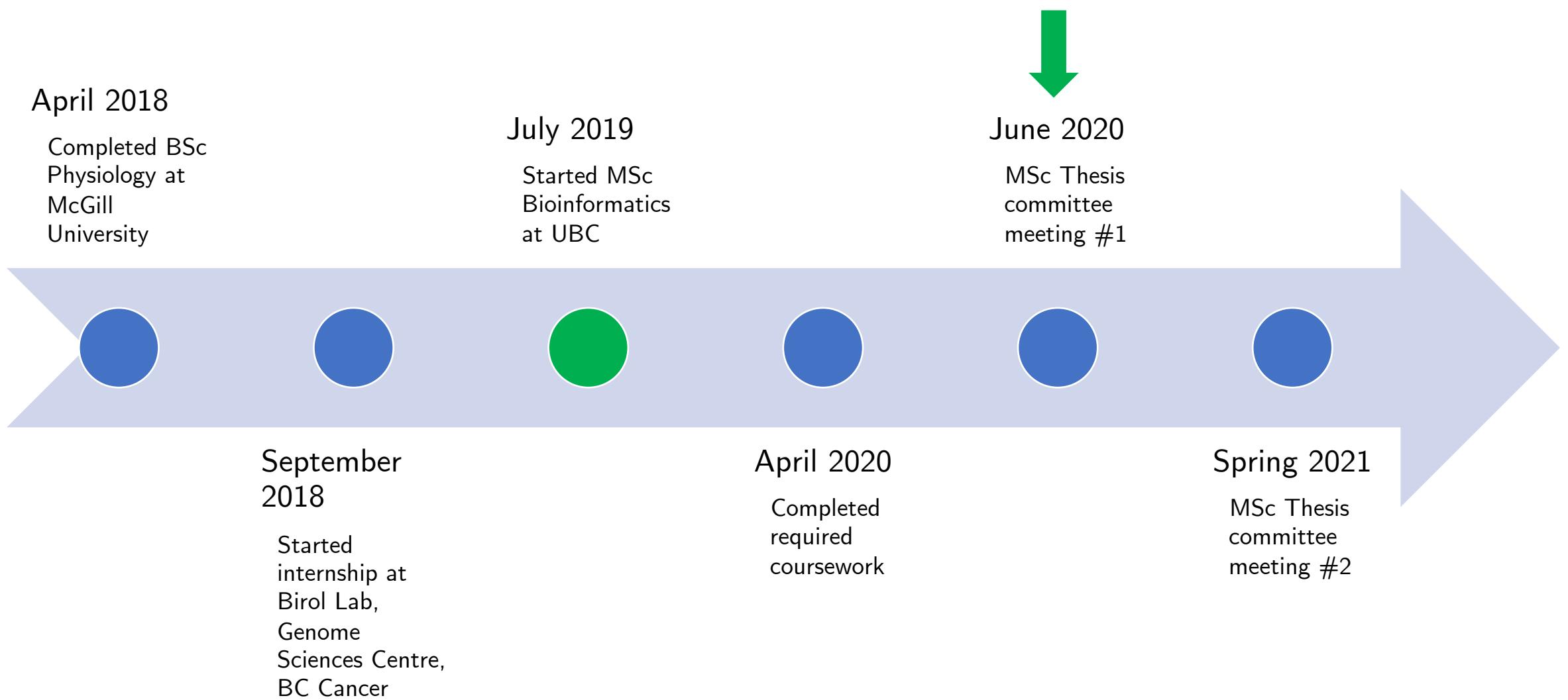
# Courses



THE UNIVERSITY  
OF BRITISH COLUMBIA

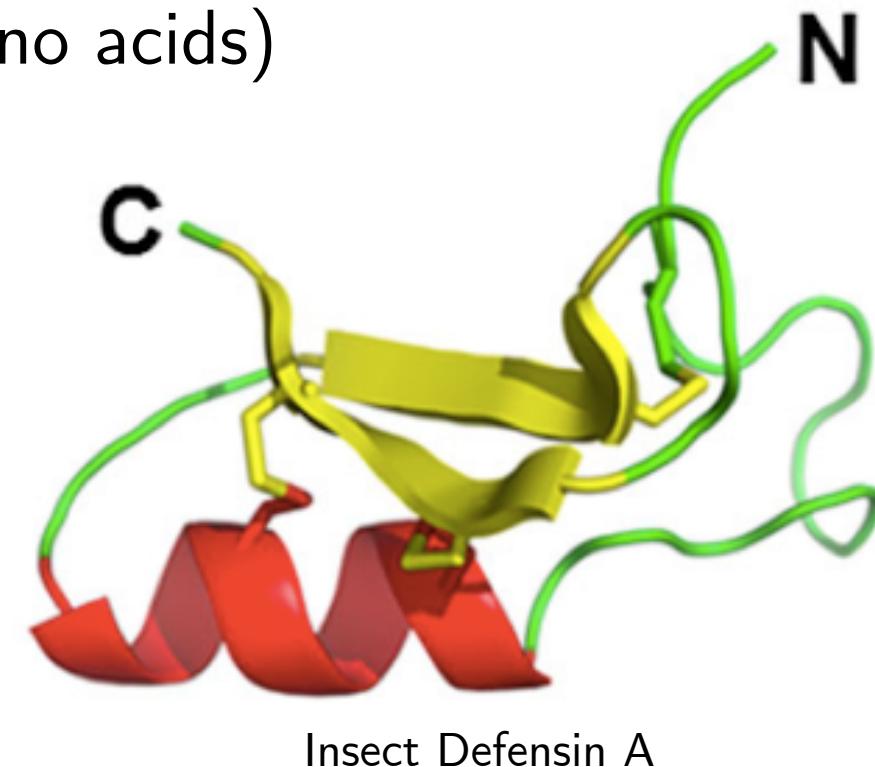
- Term 1
  - ✓ BIOF 501A: Special Topics in Bioinformatics (A)
  - ✓ STAT 545A: Exploratory Data Analysis (A+)
  - ✓ MICB 405: Introduction to Bioinformatics (A+)
- Term 2
  - ✓ BIOF 520: Problem-Based Learning in Bioinformatics (A)
  - ✓ STAT 540: Statistical Methods for High Dimensional Biology (A+)
  - ✓ STAT 547M: Topics in Statistics (A+)
  - ✓ MEDG 505: Genome Analysis (A)

# Timeline

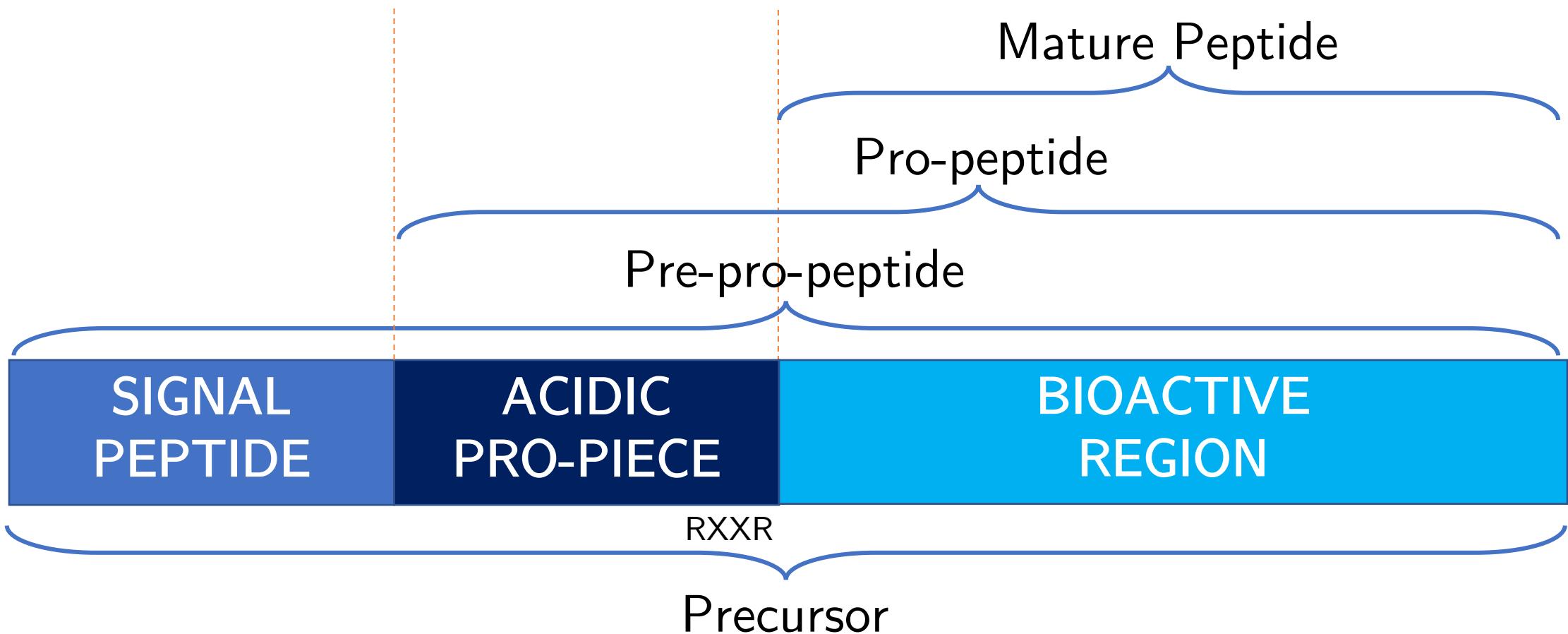


# Antimicrobial Peptides (AMPs)

- Short peptide sequences (5 to 50 amino acids)
- Often positively charged
- Amphipathic
- Produced by all life forms
- Part of the innate immune system
- Mechanisms of action:
  1. Direct interaction
  2. Modulation of host immunity



# Antimicrobial Peptides (AMPs)



# Motivation

- Rise of antibiotic resistance creates a problem that requires a novel method to fight pathogens

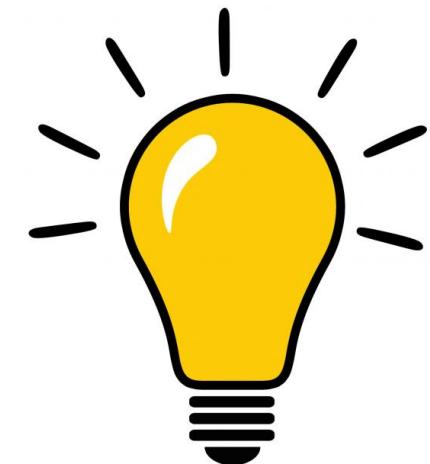
## PROBLEM

Antibiotic  
Resistance



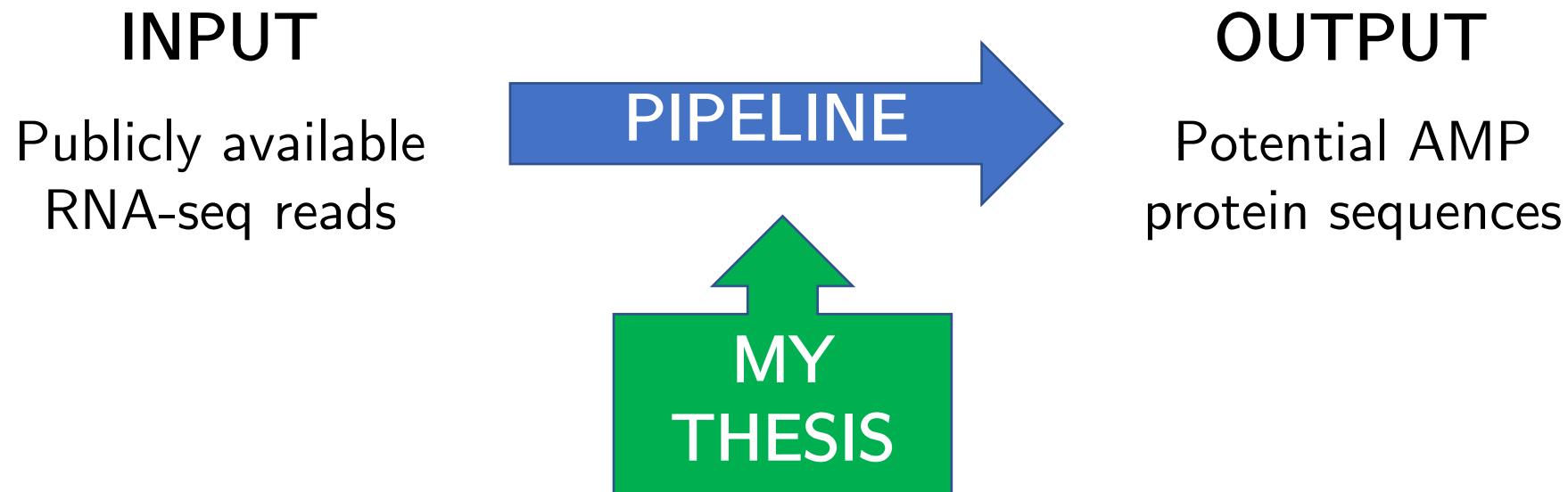
## SOLUTION

Antimicrobial  
Peptides?

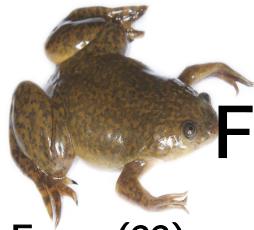


# Objective

- To develop **and** execute an AMP discovery pipeline to mine for AMP precursor sequences in publicly available genomic resources



# Datasets



## Frogs & Toads (38)



### Frogs (33)

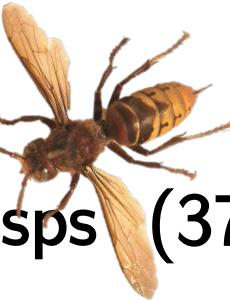
<i>A. femoralis</i>	<i>P. adspersus</i>	<i>R. sirensis</i>
<i>A. mantzorum</i>	<i>P. amboli</i>	<i>R. sylvatica</i>
<i>A. petersi</i>	<i>P. megacephalus</i>	<i>R. temporaria</i>
<i>C. alboguttata</i>	<i>P. microps</i>	<i>S. ruber</i>
<i>D. auratus</i>	<i>P. nigromaculatus</i>	<i>X. allofraseri</i>
<i>D. leucomelas</i>	<i>P. toftae</i>	<i>X. borealis</i>
<i>D. tinctorius</i>	<i>Q. boulengeri</i>	<i>X. laevis</i>
<i>H. pugnax</i>	<i>R. catesbeiana</i>	<i>X. largeni</i>
<i>L. verreauxii</i>	<i>R. dennysi</i>	<i>X. tropicalis</i>
<i>O. margaretae</i>	<i>R. imitator</i>	
<i>O. sylvatica</i>	<i>R. omeimontis</i>	
<i>O. tormota</i>	<i>R. pipiens</i>	

### Toads (5)

<i>A. minuta</i>
<i>B. gargarizans</i>
<i>L. Boringii</i>
<i>M. sangzhiensis</i>
<i>O. rhodostigmatus</i>



## Ants, Bees, & Wasps (37)



### Ants (8)

<i>A. echinatior</i>
<i>C. castaneus</i>
<i>C. obscurior</i>
<i>M. gulosa</i>
<i>O. monticola</i>
<i>P. barbatus</i>
<i>T. bicarinatum</i>
<i>T. rugulatus</i>

### Bees (5)

<i>A. cerana</i>
<i>A. mellifera</i>
<i>B. ardens</i>
<i>B. consobrinus</i>
<i>B. ussurensis</i>

### Wasps (24)

<i>A. compressa</i>	<i>P. snelleni</i>
<i>A. flavomarginatum</i>	<i>P. turionellae</i>
<i>B. nigricans</i>	<i>P. varia</i>
<i>C. vestalis</i>	<i>P. vindemmiae</i>
<i>D. collaris</i>	<i>S. deformae</i>
<i>D. longicaudata</i>	<i>S. kj8906</i>
<i>M. demolitor</i>	<i>T. sarcophagae</i>
<i>N. giraulti</i>	<i>U. rufipes</i>
<i>N. vitripennis</i>	<i>V. analis</i>
<i>N. vitripennis x N. giraulti</i>	<i>V. crabro</i>
<i>O. decorates</i>	<i>V. dybowskii</i>
<i>P. rothneyi</i>	<i>V. similis</i>

# Why Amphibians and Insects?



Source: Darcy Sutherland, PeptAID AIM 1 Science Meeting, dkfindout.com

# Datasets



## Frogs & Toads (38)



### Frogs (33)

*A. femoralis*      *P. adspersus*

*A. mantzorum*

*A. petersi*      *P. megacephalus*

***C. alboguttata***

*D. auratus*

*D. leucomelas*

*D. tinctorius*

*H. pugnax*

*L. verreauxii*

***O. margaretae***

*O. sylvatica*

*O. tormota*

*R. sirensis*

*R. sylvatica*

***R. temporaria***

*S. ruber*

***P. nigromaculatus*** *X. allofraseri*

*P. toftae*

***Q. boulengeri***

*R. catesbeiana*

*R. dennysi*

*R. imitator*

*R. omeimontis*

***R. pipiens***

### Toads (5)

*A. minuta*

*B. gargarizans*

*L. boringii*

*M. sangzhiensis*

*O. rhodostigmatus*



## Ants, Bees, & Wasps (37)

### Ants (8)

***A. echinatior***

***C. castaneus***

***C. obscurior***

*M. gulosa*

***O. monticola***

***P. barbatus***

*T. bicarinatum*

***T. rugulatus***

### Bees (5)

***A. cerana***

***A. mellifera***

*B. ardens*

*B. consobrinus*

*B. ussurensis*

*T. rugulatus*

*T. bicarinatum*

*T. rugulatus*

*N. giraulti*

***N. vitripennis***

*N. vitripennis* x

*N. giraulti*

*O. decorates*

*P. rothneyi*

### Wasps (24)

*A. compressa*

*A. turionellae*

*B. nigricans*

*C. vestalis*

*D. collaris*

*D. longicaudata*

*M. demolitor*

*N. giraulti*

***N. vitripennis***

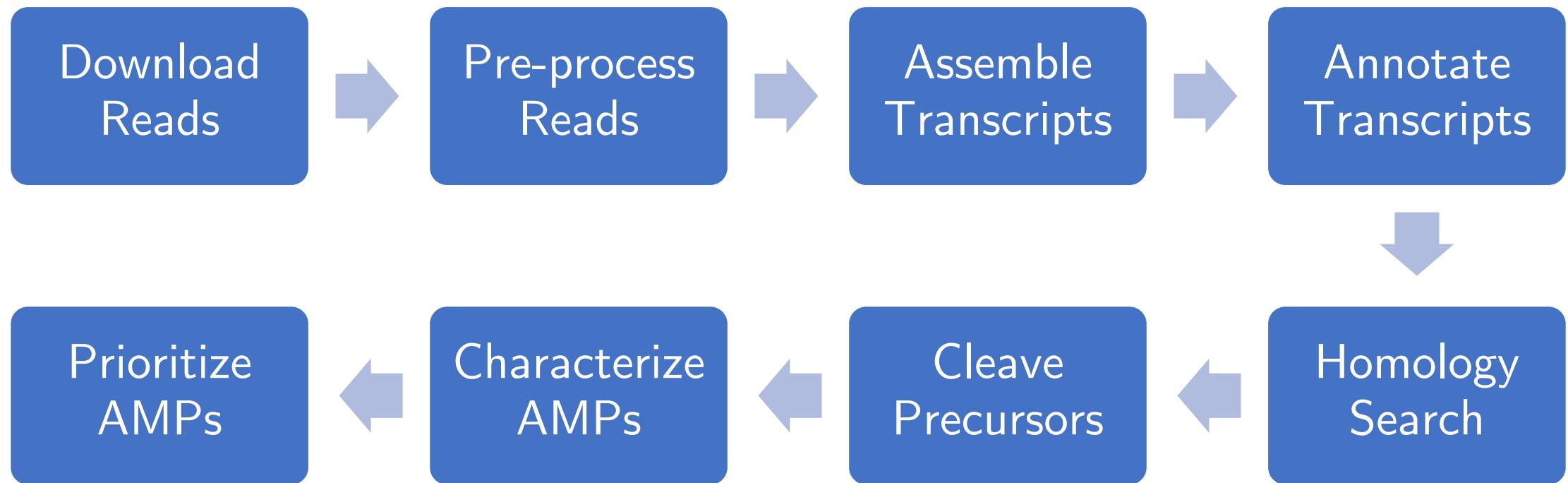
*V. crabro*

*V. dybowskii*

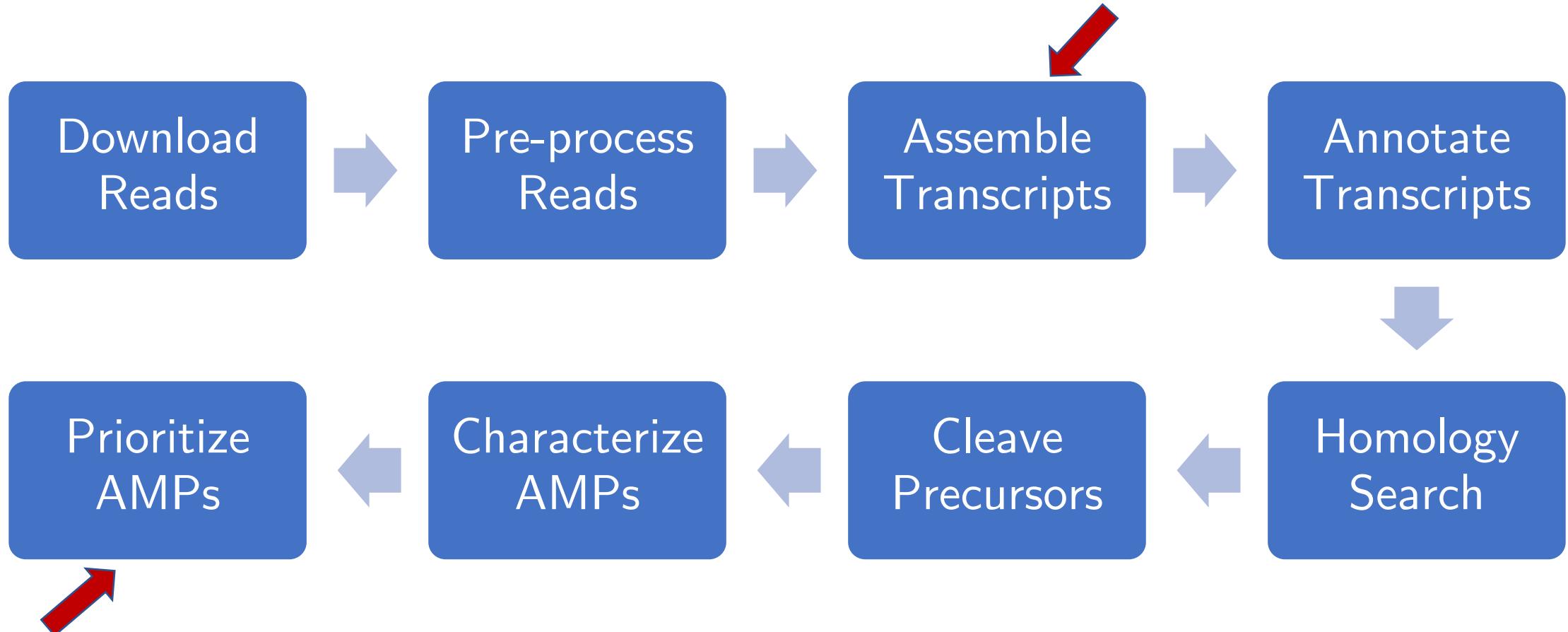
*V. analis*

*V. similiia*

# Methods Overview



# Methods Overview



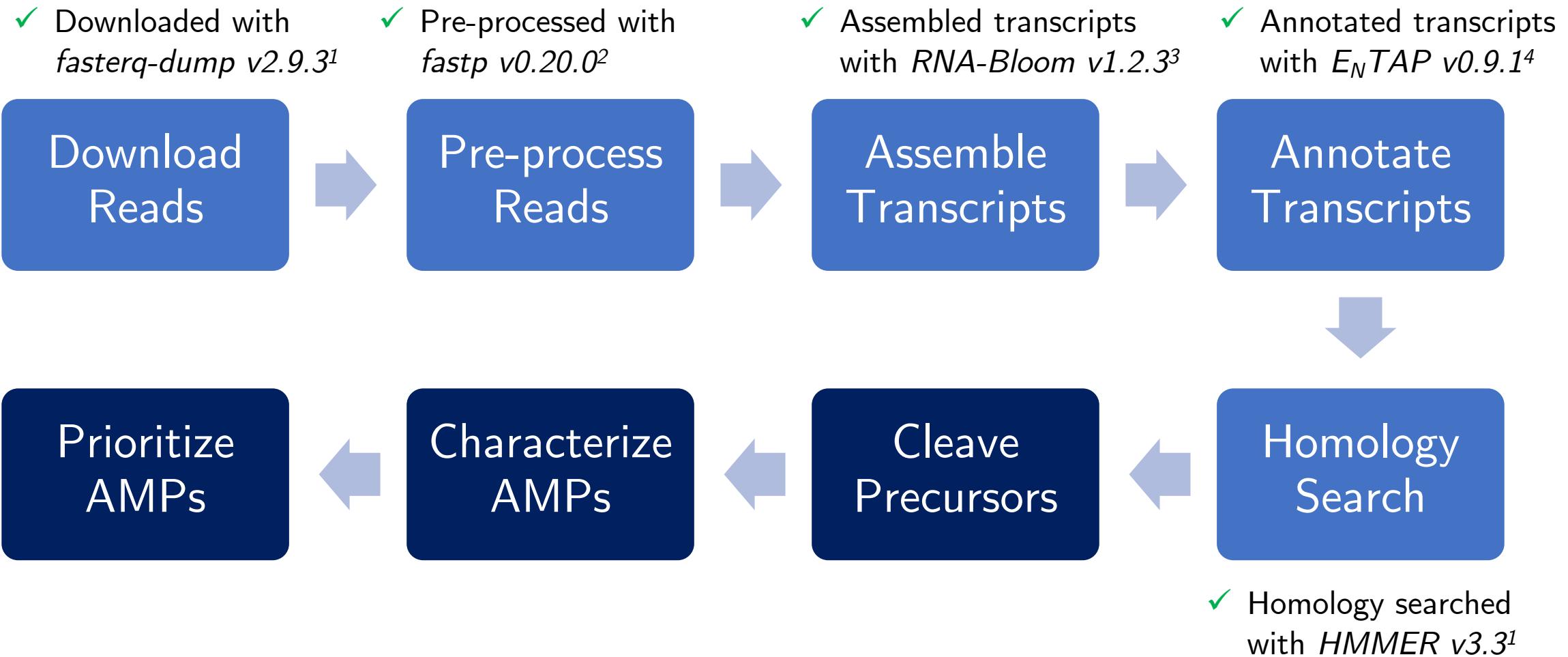
RNA-Bloom

<https://github.com/bcgsc/RNA-Bloom>

AMPlify

<https://github.com/bcgsc/amplify>

# Progress





# Pre-processing: fastp

Chen, S., Zhou, Y., Chen, Y. & Gu, J. **fastp: an ultra-fast all-in-one FASTQ preprocessor**. Bioinformatics 34, i884–i890 (2018).

- Quality control - FASTQC
- Adapter trimming - Cutadapt
- Quality filtering - Trimmomatic



fastp **replaces** the need  
to use *three separate tools*

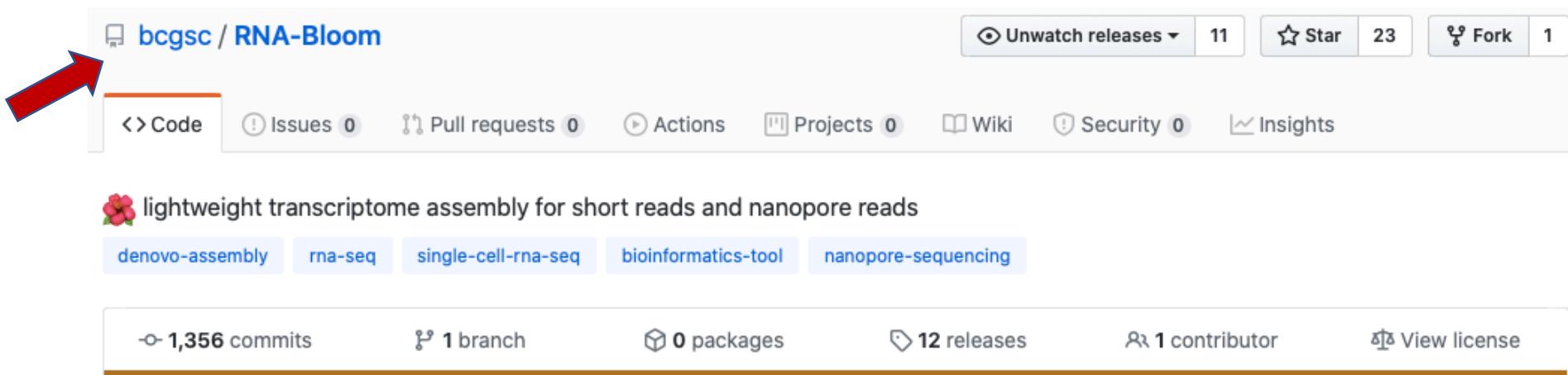


# Assembly: RNA-Bloom



Nip, K. M. et al. RNA-Bloom provides lightweight reference-free transcriptome assembly for single cells. Genome Research (under revision).

- *De novo* transcriptome assembly with single and paired-end reads
- Reference-guided assembly if reference or draft transcriptome available



bcgsc / RNA-Bloom

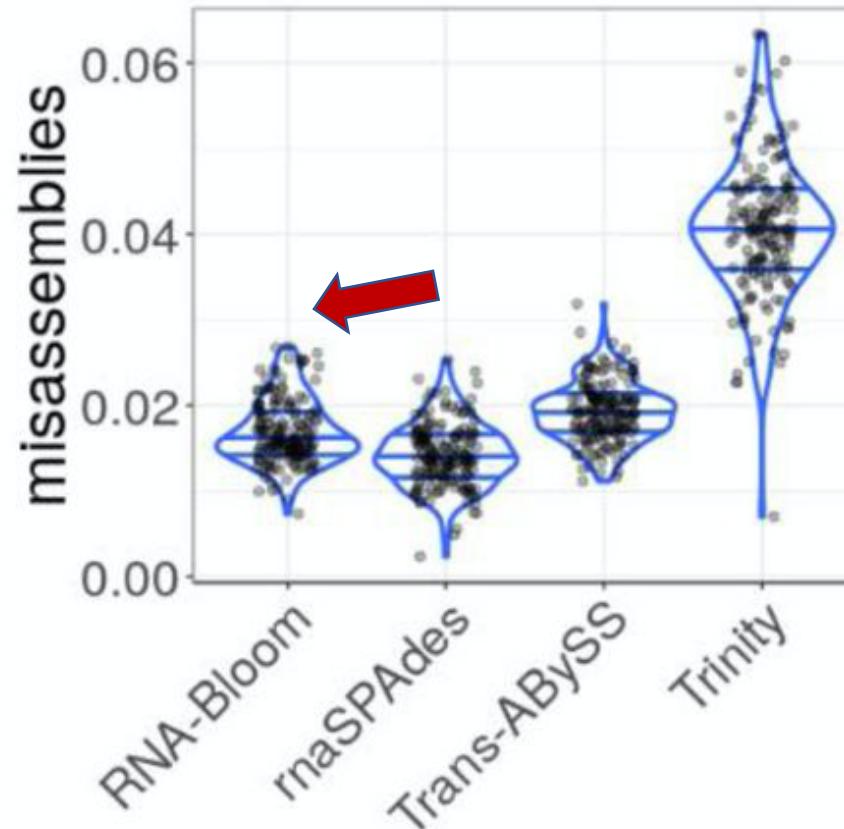
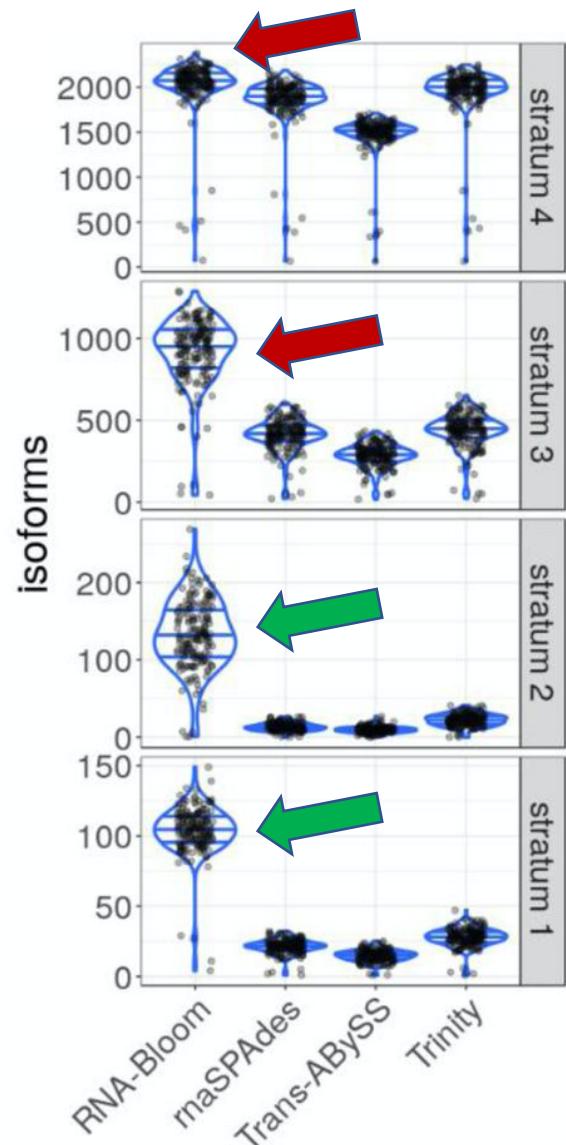
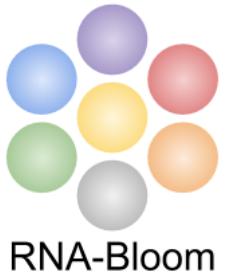
Code Issues 0 Pull requests 0 Actions Projects 0 Wiki Security 0 Insights

lightweight transcriptome assembly for short reads and nanopore reads

denovo-assembly rna-seq single-cell-rna-seq bioinformatics-tool nanopore-sequencing

1,356 commits 1 branch 0 packages 12 releases 1 contributor View license

# Assembly: RNA-Bloom



Source: Nip, K. M. et al. RNA-Bloom provides lightweight reference-free transcriptome assembly for single cells. *Genome Research* (under revision).

# Annotation: E<sub>N</sub>TAP

Hart, A. J. et al. E<sub>N</sub>TAP: Bringing faster and smarter functional annotation to non-model eukaryotic transcriptomes. Mol. Ecol. Resour. 20, 591–604 (2020).

- E<sub>N</sub>TAP: Eukaryotic Non-Model Transcriptome Annotation Pipeline
- Functional annotation with ‘frame selection’
- Alignment to orthology, domain, and protein databases

# Annotation: ENTAP

Metric	BLAST2GO PRO	BLAST2GO Basic	TRINOTATE	ENTAP	ANNOCRIPT	DAMMIT
Open source/free software		‡	‡	‡	‡	‡
Command line integration	‡		‡	‡	‡	‡
Filtering assembly via short-read alignment (expression)	†				‡	
Frame selection	†		‡	‡	‡	‡
Custom database selection and indexing	‡		†	‡		†
Fast and sensitive NCBI BLAST alternative				‡	†	
Selection of optimal hit From several databases				‡	†	
Selection of optimal hit based on informativeness	†				‡	
Contaminant identification and filtering	‡				‡	
Orthologous gene family assignment	‡		†	‡		
Protein domain (CDD/InterProScan)	‡	‡	‡	‡	‡	‡
Gene ontology term and pathway assignment Sourced from protein alignments	‡	‡	‡		‡	
Gene ontology term and pathway assignment Sourced from orthologous genes	‡			‡		‡
Provides graphical user interface for annotation process	‡	‡				

Source: Hart, A. J. et al. ENTAP: Bringing faster and smarter functional annotation to non-model eukaryotic transcriptomes. Mol. Ecol. Resour. 20, 591–604 (2020).



# Homology Search: HMMER

Johnson, L. S., Eddy, S. R. & Portugaly, E. Hidden Markov model speed heuristic and iterative HMM search procedure. BMC Bioinformatics 11, 431 (2010).

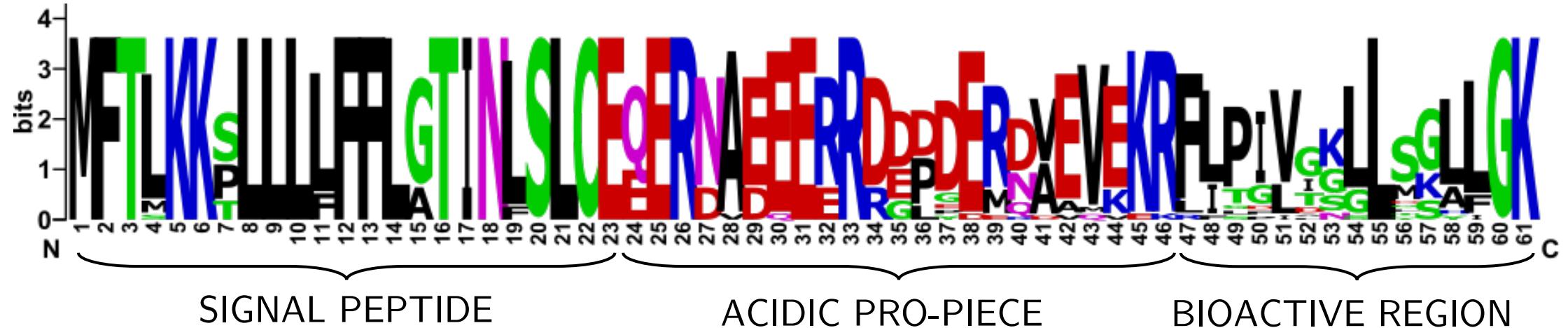
- Using antimicrobial precursor sequences from NCBI Protein database (including GenPept, UniProt, etc.)
- 2,792 amphibian precursors; 303 insect precursors



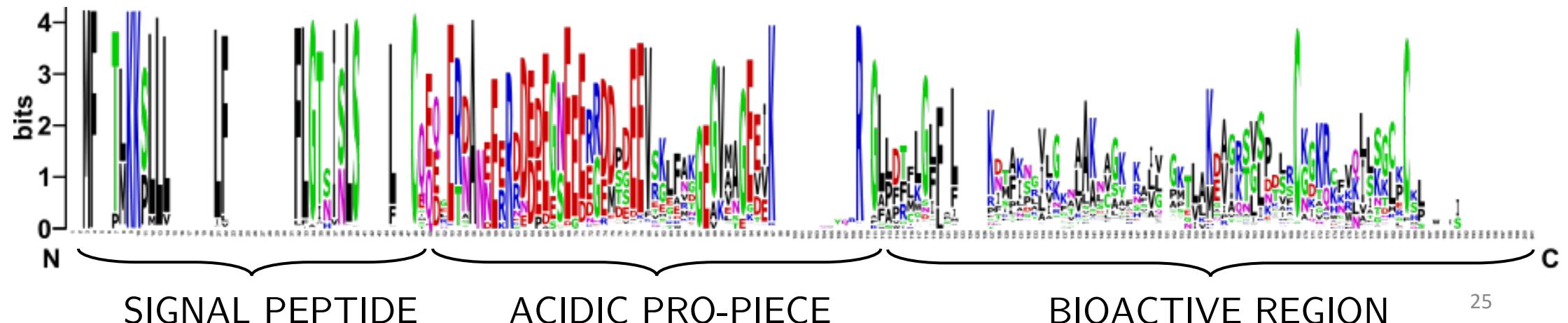


# Homology Search: HMMER

Single AMP family:

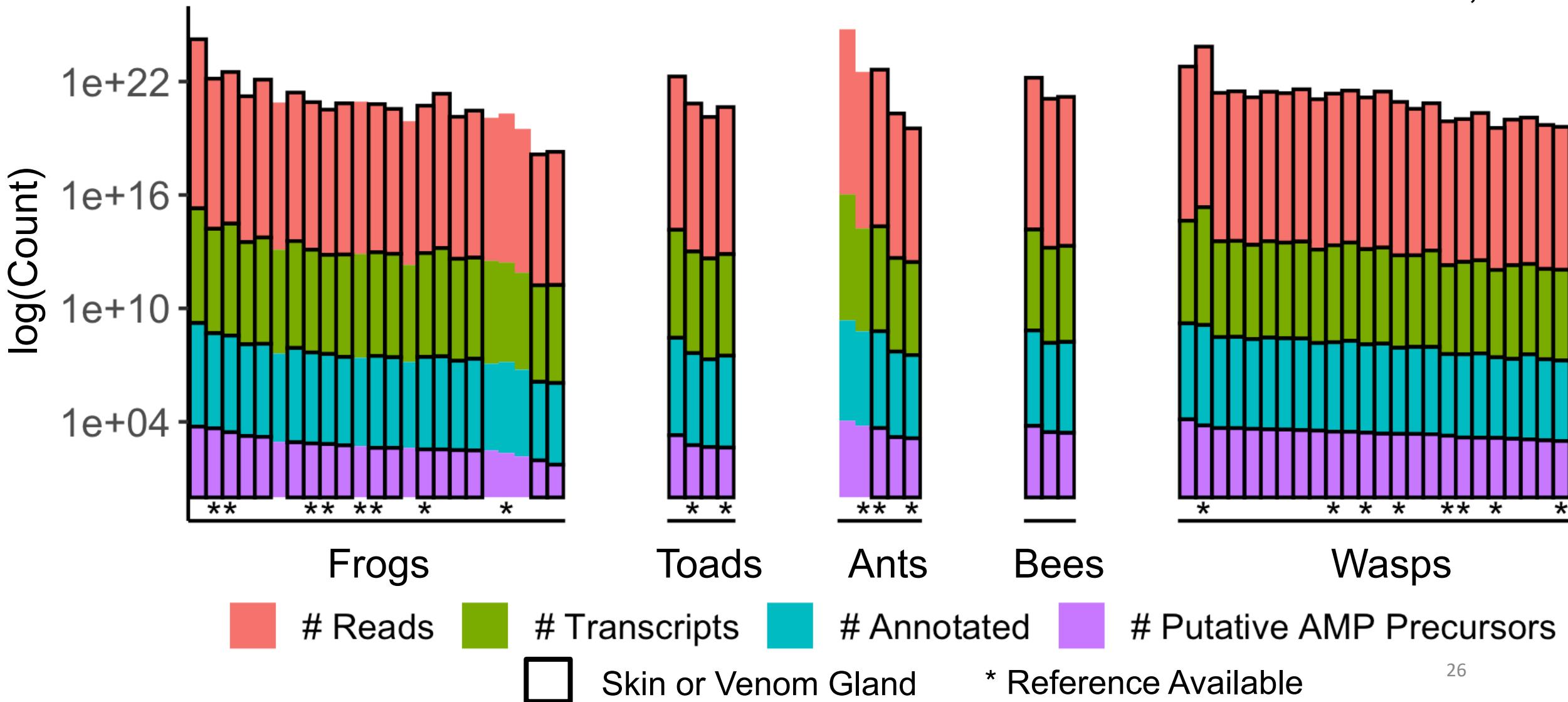


General AMP structure:

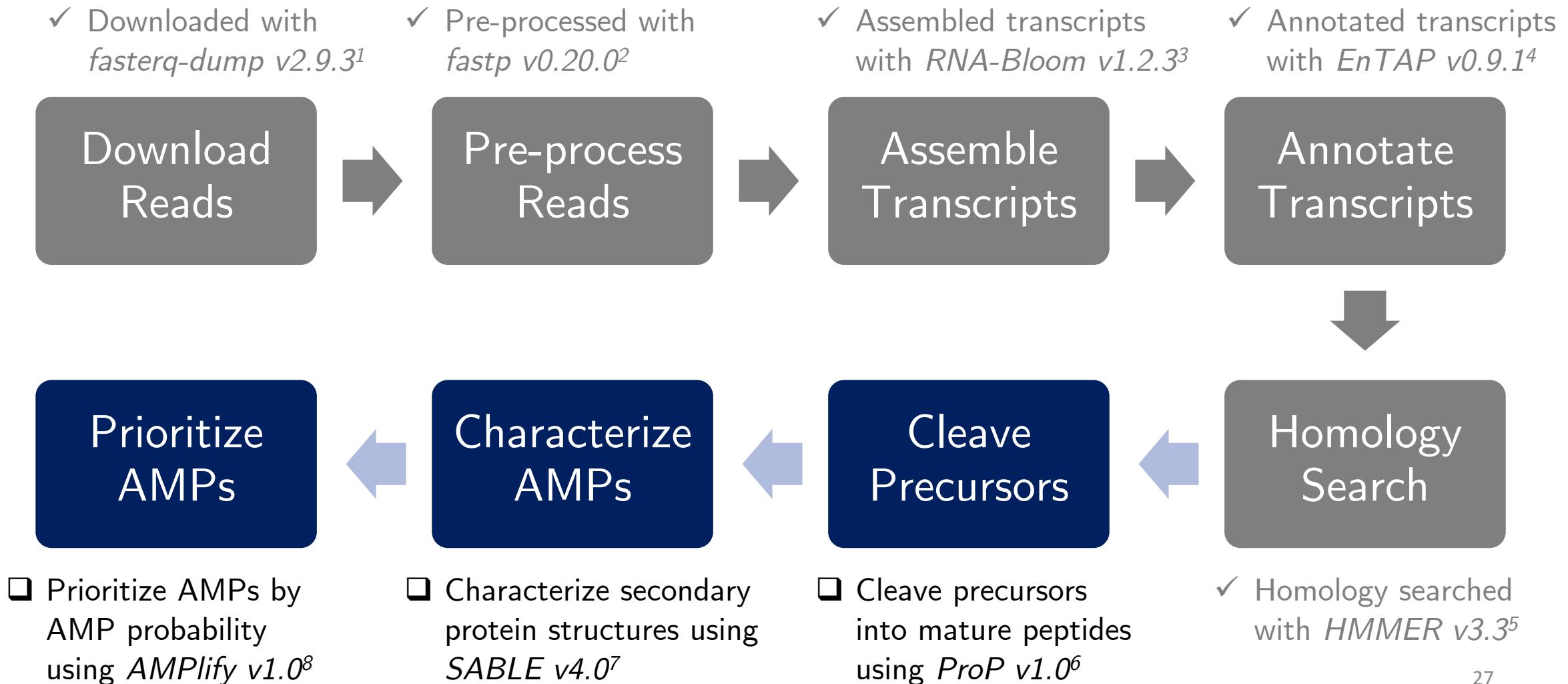


# Preliminary Results

Total # Putative AMP Precursors: 140,054

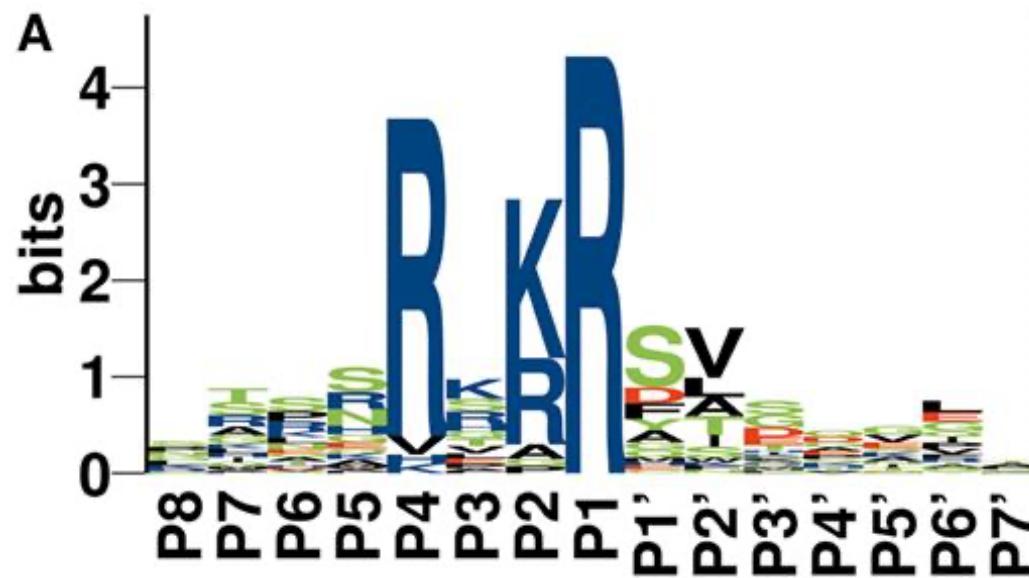


# Future Work



# Precursor Cleavage: ProP

Duckert, P., Brunak, S. & Blom, N. Prediction of proprotein convertase cleavage sites. Protein Eng. Des. Sel. 17, 107–112 (2004).



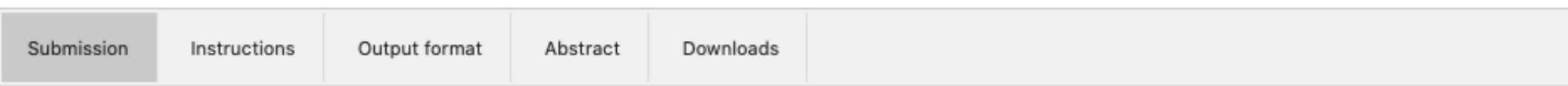
# Precursor Cleavage: ProP

## ProP - 1.0

### Arginine and lysine propeptide cleavage sites in eukaryotic protein sequences

The ProP 1.0 server predicts arginine and lysine propeptide cleavage sites in eukaryotic protein sequences using an ensemble of neural networks. Furin-specific prediction is the default. It is also possible to perform a general proprotein convertase (PC) prediction.

For convenience, this server is integrated with the [SignalP-3.0](#) server predicting the presence and location of signal peptide cleavage sites.



>sp|P10891.2|DEFI\_PROTE RecName: Full=Phormicin; AltName: Full=Insect defensin-A/B; Flags: Precursor

MKFFMFVVTFCLAVCFVSQSLAIPADAANDAHFVDGVQALKEIE  
PELHGRYKRATCDLLSGTGINHSACAAHCLLRGNRGGYCNGKGVC  
VCRN

Source: Duckert, P., Brunak, S. & Blom, N. Prediction of proprotein convertase cleavage sites. Protein Eng. Des. Sel. 17, 107–112 (2004).

# Precursor Cleavage: ProP

Sequence: MKFFMVFVVTFCLAVCFVSQSLAIPADAANDAHFVDGVQALKEIEPELHGRYKRATCDILLSGTGINHSACAAHCLLRGNR

Signal peptide cleavage site: MKFFMVFVVTFCLAVCFVSQSLAIPADAANDAHFVDGVQALKEIEPELHGRYKRATCDILLSGTGINHSACAAHCLLRGNR

Propeptide cleavage site: MKFFMVFVVTFCLAVCFVSQSLAIPADAANDAHFVDGVQALKEIEPELHGRYKRATCDILLSGTGINHSACAAHCLLRGNR

Signal peptide: MKFFMVFVVTFCLAVCFVSQSLAIPADAANDAHFVDGVQALKEIEPELHGRYKRATCDILLSGTGINHSACAAHCLLRGNR

acidic pro-piece: MKFFMVFVVTFCLAVCFVSQSLAIPADAANDAHFVDGVQALKEIEPELHGRYKRATCDILLSGTGINHSACAAHCLLRGNR

Signal peptide cleavage site predicted: between pos. 23 and 24: SLA-IP

Propeptide cleavage sites predicted: Arg(R)/Lys(K): 1

>sp|P10891.2|DEFI\_PROTE RecName: Full=Phormicin; AltName: Full=Insect defensin A/B; Flags: Precursor

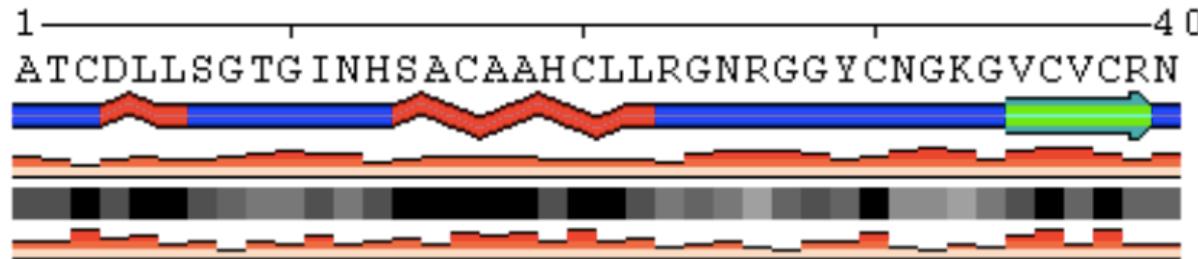
MKFFMVFVVTFCLAVCFVSQSLAIPADAANDAHFVDGVQALKEIEPELHGRYKR  
ATCDLLSGTGINHSACAAHCLLRGNRGGYCNGKGVCVRN

Source: Duckert, P., Brunak, S. & Blom, N. Prediction of proprotein convertase cleavage sites. Protein Eng. Des. Sel. 17, 107–112 (2004). 30

# Characterization: SABLE

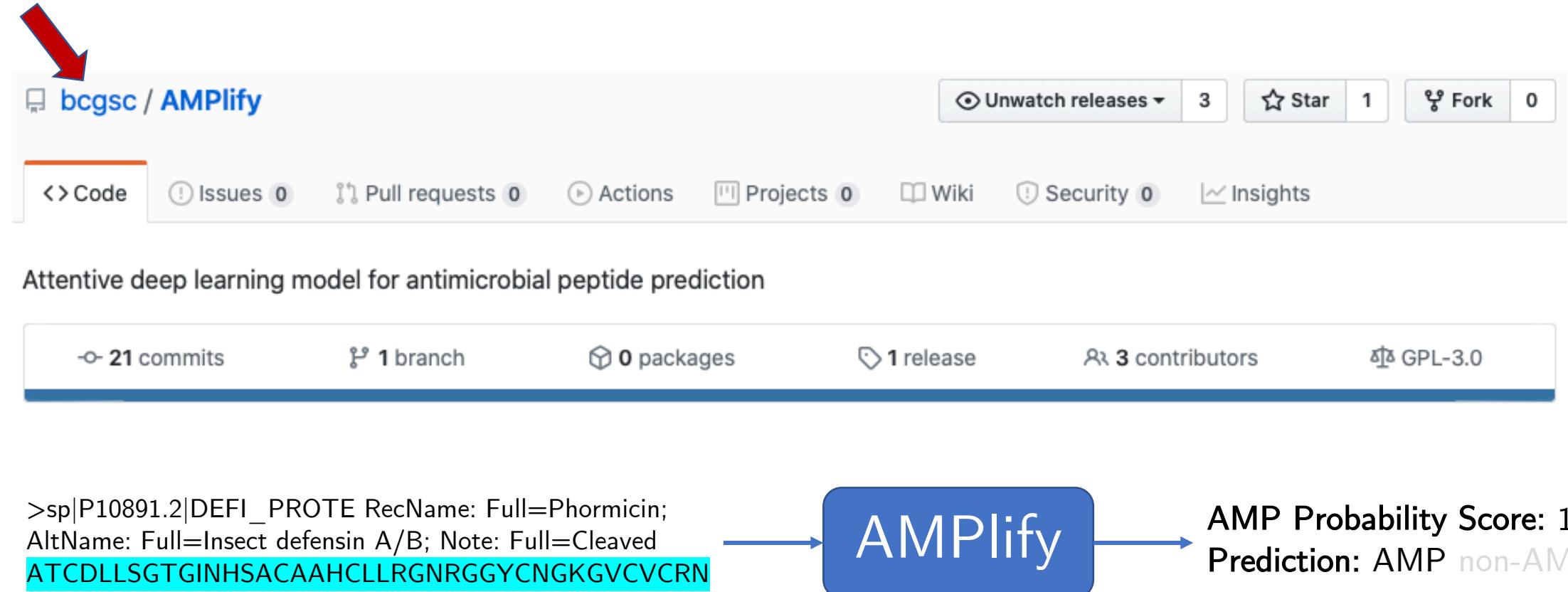
Adamczak, R., Porollo, A. & Meller, J. Combining prediction of secondary structure and solvent accessibility in proteins. Proteins 59, 467–475 (2005).

>sp|P10891.2|DEFI\_PROTE RecName: Full=Phormicin;  
AltName: Full=Insect defensin A/B; Note: Full=Cleaved  
**ATCDLLSGTGINHSACAAHCLLRGNRGGYCNGKGVCVRN**



Legend	Description
1 —————	Amino acid residue numeration
	Protein secondary structure
	H-alpha and other helices
	E-beta-strand or bridge
	Relative solvent accessibility (RSA)
0 1 2 3 4 5 6 7 8 9	0-completely buried (0-9% RSA), 9-fully exposed (90-100% RSA)
	Confidence level of prediction
0 1 2 3 4 5 6 7 8 9	0-the lowest level, 9-the highest level

# Prioritization: AMPlify



Source: Li, C. et al. AMPlify: Attentive Deep Learning Model for Discovery of Novel Antimicrobial Peptides Effective against WHO Priority Pathogens. Nature Methods (submitted).

# End Product

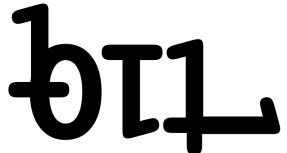
- ✓ AMP Discovery Pipeline software package that runs all the tools from beginning to end
- ✓ Resulting from the pipeline, candidate AMPs for downstream analysis and *in vitro* bioactivity testing against various microbes

# References

1. Leinonen, R., Sugawara, H., Shumway, M. & International Nucleotide Sequence Database Collaboration. The sequence read archive. *Nucleic Acids Res.* **39**, D19–21 (2011).
2. Chen, S., Zhou, Y., Chen, Y. & Gu, J. fastp: an ultra-fast all-in-one FASTQ preprocessor. *Bioinformatics* **34**, i884–i890 (2018).
3. Nip, K. M. et al. RNA-Bloom provides lightweight reference-free transcriptome assembly for single cells. *Genome Research* (under revision). 
4. Hart, A. J. et al. E<sub>N</sub>TAP: Bringing faster and smarter functional annotation to non-model eukaryotic transcriptomes. *Mol. Ecol. Resour.* **20**, 591–604 (2020).
5. Johnson, L. S., Eddy, S. R. & Portugaly, E. Hidden Markov model speed heuristic and iterative HMM search procedure. *BMC Bioinformatics* **11**, 431 (2010).
6. Duckert, P., Brunak, S. & Blom, N. Prediction of proprotein convertase cleavage sites. *Protein Eng. Des. Sel.* **17**, 107–112 (2004).
7. Adamczak, R., Porollo, A. & Meller, J. Combining prediction of secondary structure and solvent accessibility in proteins. *Proteins* **59**, 467–475 (2005).
8. Li, C. et al. AMPlify: Attentive Deep Learning Model for Discovery of Novel Antimicrobial Peptides Effective against WHO Priority Pathogens. *Nature Methods* (submitted). 

# Acknowledgements

Birol Lab



- Ka Ming Nip
- Kristina Gagalova
- Darcy Sutherland
- Chenkai Li
- René Warren
- Inanc Birol

PeptAID Collaborators

- Caren Helbing

PeptAID Funding



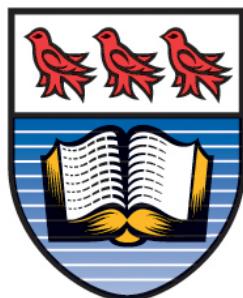
**Genome**  
BritishColumbia



**Genome**Canada



**Investment  
Agriculture  
Foundation  
of British Columbia**



**University  
of Victoria**