National Institute of Allergy and Infectious Diseases

**2019 Health Innovations Conference**

# Using deep learning to improve antimicrobial peptide recognition
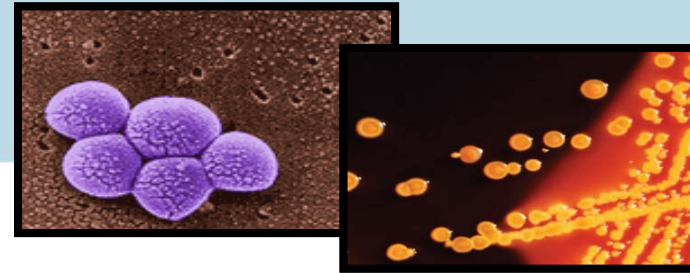
**Tuesday, 19 March 2019**

NIAID

National Institute of Allergy and Infectious Diseases

NIH

**Daniel Veltri, Ph.D.**

# Motivation



Methicillin-resistant *Staphylococcus aureus* (left)
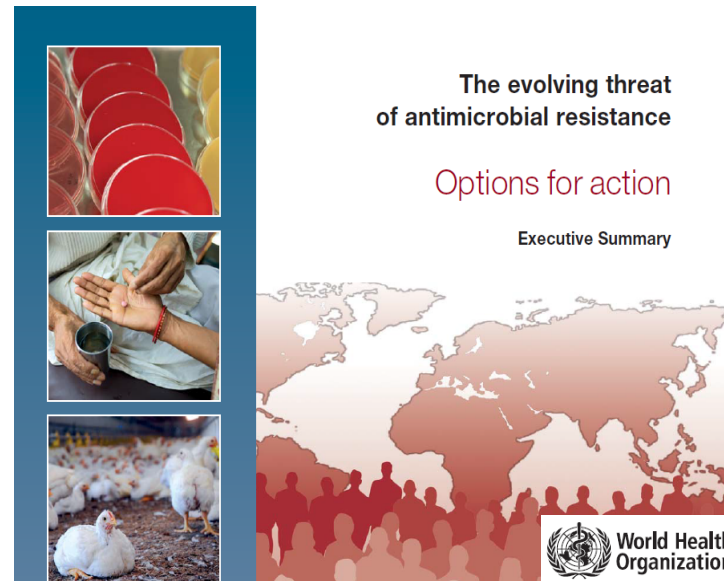Carbapenem-resistant Enterobacteriaceae (right)

- Reports of antibiotic resistance have *increased* despite a *slowdown* in new antibiotics coming to market in recent decades,

- The U.S. Center for Disease Control reports over 2 million infections and 23,000 deaths each year due to antibiotic-resistant bacteria and fungi in the U.S. [1],

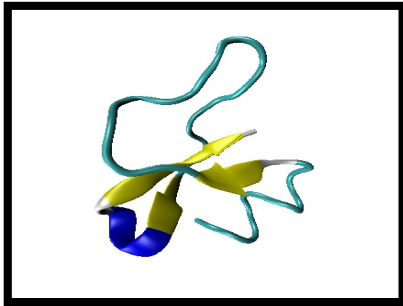- The WHO has put out numerous reports warning of the risks of resistant bacteria to hospitals around the world.
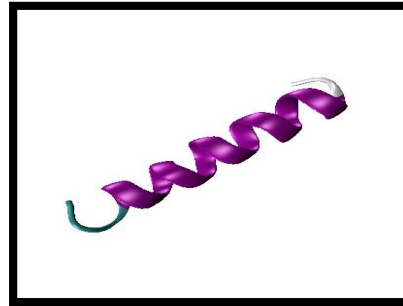
[1] CDC: https://www.cdc.gov/drugresistance/about.html



The evolving threat of antimicrobial resistance

Options for action

Executive Summary

World Health Organization

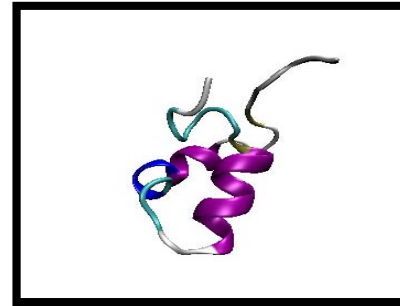National Institute of Allergy and Infectious Diseases

NIAID

# Antimicrobial Peptides (AMPs)
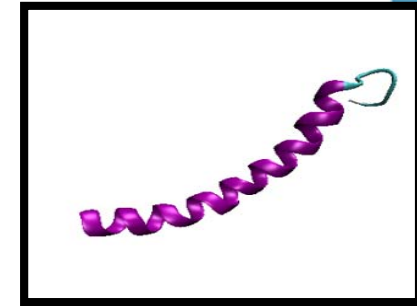


**Beta Defensin 1**
*Homo sapiens*
PDB: 1IJV

**Magainin 2**
*Xenopus laevis*
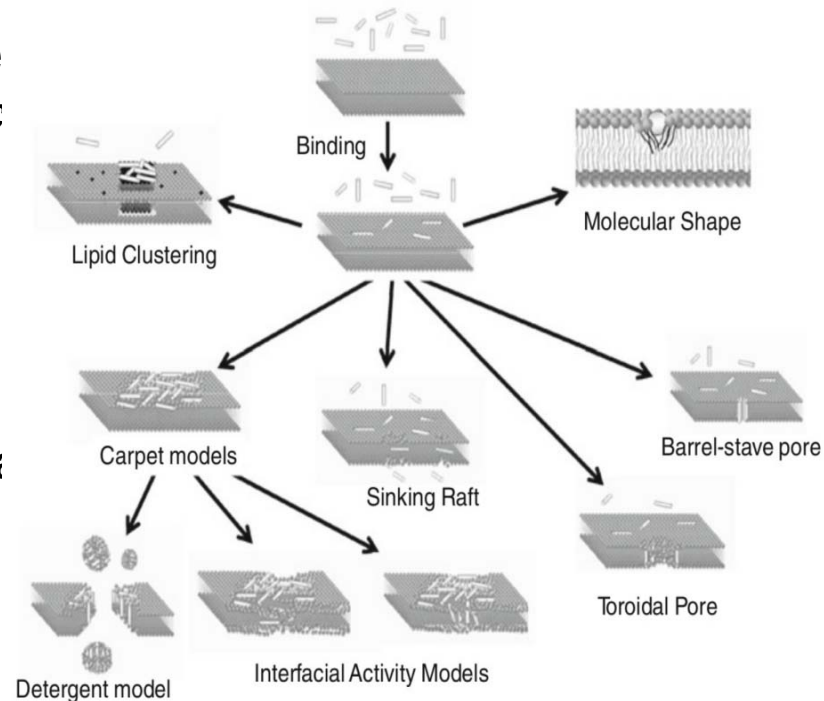PDB: 2MAG

**Aurelin**
*Aurelia aurita*
PDB: 2LG4

**Cathelicidin LL-37**
*Homo sapiens*
PDB: 2K6O

- One promising area for new antibiotic research has been natural AMPs- short peptides with innate antibacterial activity found across all phyla,

- To date, efforts to design and/or modify AMPs have had limited success in delivering new drugs to market.

National Institute of
Allergy and
Infectious Diseases

NIH

NIAID

3

# AMPs are Complicated!

- Amino acid (AA) physicochemical properties are important for AMP activity (*charge, hydrophobic etc.*),

- AMPs are highly diverse, both in sequence and killing mechanism,

- We still do not know *exactly* how physicochemical properties relate to AMP activity- knowledge needed to *guide AMP design.*

National Institute of Allergy and Infectious Diseases

Some proposed AMP attack mechanisms

4

# Prior AMP Classification Work

- Most work to date has focused on AMP recognition- taking query peptide sequences and assigning *AMP* or *non-AMP* labels,

- Top techniques report accuracies in the high 80 to mid 90% range,

- Approaches often pair physicochemical properties with sliding window averages or machine learning algorithms like artificial neural networks (ANN), support vector machines (SVM), etc.,

- A major issue in the field is that *few groups make their code or complete data sets available*. This makes it difficult to perform reliable comparisons as a "gold standard" benchmark data set is not currently available.

National Institute of
Allergy and
Infectious Diseases

NIH

5

# Prior AMP Classification Performance

| Group | Algorithm | Performance |
|-------|-----------|-------------|
| Ch... | | |

Matthew's Correlation Coefficient (MCC):

$$MCC = \frac{(TP \text{ x } TN) - (FN \text{ x } FP)}{\sqrt{(TP + FN)\text{x}(TN + FP)\text{x}(TP + FP)\text{x}(TN + FN)}}$$

*MCC values range from -1 to 1, with 1 denoting perfect classification performance.*

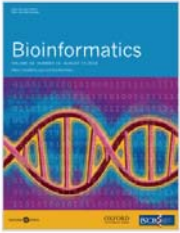TP = True Positive, TN = True Negative, FP = False Positive, FN = False Negative

| | | |
|-------|-----------|-------------|
| Xiao et al. (2013) | Fuzzy K-Nearest Neighbor | 0.84 |
| Meher et al. (2017) | SVM | 0.84 |

National Institute of Allergy and Infectious Diseases

† "Random Forest" is trademarked and licensed to Salford Systems (San Diego, CA)

# Using Deep Learning for AMP Classification

# Deep Neural Networks (DNN) in the News...



Go board image from Wikimedia Foundation

## Deep Learning Packages
### *So many flavors to choose from...*

# Deep Neural Networks Have Multiple Layers



Source: Chollet and Allaire "Deep Learning with R" pp.9, 2018.

# Our Model Architecture

# Convolutional Layers



*Figure from: towardsdatascience.com*

# Pooling Layer



*Figure from: towardsdatascience.com*

# Long Short-Term Memory (LSTM)



*Direction of Reading →*

**Ignore!**

... TCCGCGATCGTTCGGGTGGCCTTTAATATTATGTGCGCGTTAGCTGGTCACGCG

**Recognize Pattern!**

**Original LSTM Paper:** Hochreiter and Schmidhuber (1997) Long short-term memory.

**Figure:** *deeplearning.net*

14

# Data Set Construction

- AMPs were taken from the Antimicrobial Peptide Database vr3 (aps.unmc.edu/AP). Removed any <10 AA in length or sharing ≥90% sequence identity,

- Non-AMPs taken from UniProt using keyword filtering. Removed any <10 AA in length or sharing ≥40% sequence identity,

- Randomly selected even number of AMPs and Non-AMPs for each partition: 712 Training, 354 Tuning, and 712 Testing.

# Model Training and Testing Performance

| Training set | Evaluation set | SENS(%) | SPEC(%) | ACC(%) | MCC | auROC(%) |
|---|---|---|---|---|---|---|
| Train-Only | Train | 98.60 | 98.87 | 98.69 | 0.9706 | 99.87 |
| Train-Only | Tune | 95.76 | 83.85 | 87.80 | 0.7582 | 96.67 |
| Train+Tune | Train+Tune | 97.19 | 99.53 | 98.36 | 0.9674 | 99.75 |
| Train+Tune | Test | 89.89 | 92.13 | 91.01 | 0.8204 | 96.48 |
| All Data | All Data | 98.26 | 99.66 | 98.96 | 0.9793 | 99.94 |
| All Data | 10-fold CV | 88.81 (±3.53) | 94.21 (±2.68) | 91.51 (±0.89) | 0.8327 (±0.02) | 96.58 (±0.66) |

National Institute of
Allergy and
Infectious Diseases

# A Head-to-Head AMP Server Comparison

*Classification performance on our testing data set*

| Method | SENS(%) | SPEC(%) | ACC(%) | MCC | auROC(%) |
|---|---|---|---|---|---|
| AntiBP2 (imtech.res.in/raghava/antibp2) | 87.91 | 90.80 | 89.37 | 0.7876 | 89.36 |
| CAMP-ANN (camp.bicnirrh.res.in/predict) | 82.98 | 85.09 | 84.04 | 0.6809 | 84.06 |
| CAMP-DA | 87.08 | 80.76 | 83.92 | 0.6797 | 89.97 |
| CAMP-RF | **92.70** | 82.44 | 87.57 | 0.7554 | 93.63 |
| CAMP-SVM | 88.90 | 79.92 | 84.41 | 0.6910 | 90.63 |
| iAMP-2L (jci-bioinfo.cn/iAMP-2L) | 83.99 | 85.86 | 84.90 | 0.6983 | 84.90 |
| iAMPpred (cabgrid.res.in:8080/amppred) | 89.33 | 87.22 | 88.27 | 0.7656 | 94.44 |
| Our DNN | 89.89 | **92.13** | **91.01** | **0.8204** | **96.48** |

National Institute of
Allergy and
Infectious Diseases

# AMP Server Comparison ROC Curve

# Embedding Vector of Amino Acids

A 2D t-SNE [1] projection of the 128 dim. AA embedding vectors. K-means (*k=9*) used to select clusters

[1] Van der Maaten et al. J. Machine Learn. Res., 2008

# AMP Scanner vr.2 Website

Feel free to try our methods out at:

**_www.ampscanner.com_**

Data sets are available to download and contact information if you would like the code from me!

National Institute of
Allergy and
Infectious Diseases

# Building a Generative Model for AMP-like Sequences

## Guiding Exploration of Antimicrobial Peptide Space with a Deep Neural Network

Manpriya Dua
Dept of Computer Science
George Mason University
Fairfax, VA, USA
mdua@gmu.edu

Daniel Veltri
National Institute of
Allergy & Infectious Diseases
National Institutes of Health
Rockville, MD, USA
dan.veltri@gmail.com

Barney Bishop
Dept of Chemistry & Biochemistry
George Mason University
Manassas, VA, USA
bbishop1@gmu.edu

Amarda Shehu
Dept of Computer Science
George Mason University
Fairfax, VA, USA
amarda@gmu.edu

*Abstract*—Antibiotic resistance has become a serious concern, and many health organizations are sounding the alarm and the need for new drug templates. Naturally-occurring antimicrobial peptides (AMPs) have long promised to serve as such templates, as they have shown lower likelihood for bacteria to form resistance. This has motivated wet and dry laboratories to seek These peptides fall into a number of diverse sequence families (e.g. cathelicidins, defensins, cecropins, etc.), are diverse in secondary and tertiary structure, and kill their targets through various mechanisms, such as cell membrane damage, DNA interference, or signaling for adaptive immune responses [8].

National Institute of
Allergy and
Infectious Diseases

NIH

NIAID

# 4 Different Sampling Methods

- *RANDOM-* randomly select AAs and a total peptide length (L) from a population of training AMPs (baseline method),

- *GREEDY-* Perform *RANDOM*, then select (L+1) AA's to substitute with changes that improve AMP probability (use prior deep learning model to judge),

- *Metropolis Monte Carlo (MMC)-* Perform *GREEDY* but, with a small probability, accepted *worse* AA changes. Temperature (**T**) parameter decides how often we do this (higher **T** → more changes → more diverse sequences),

- **Simulated Annealing (SA-MMC)-** Similar to MMC above but starts with a high **T** to start more diverse and gradually lowers **T** over time to become greedier.

National Institute of
Allergy and
Infectious Diseases

# Distribution of *All* Generated Peptides

*The simulated annealing (SA-MCC) method performs best- it generates the most sequences predicted to be antimicrobial*

$\leftarrow$ *Predicted non-AMP* **0.5** *Predicted AMP* $\rightarrow$

# Future Directions

- Now that we can generate and evaluate AMP sequences, can we use *adversarial learning* to build improved AMP classifiers?

- More work needs to be done predicting how AMPs may work against *specific bacteria* of medical interest. Can we do better at predicting MIC, EC50 etc.?

National Institute of
Allergy and
Infectious Diseases

# Collaborators



Amarda Shehu, Ph.D.
George Mason University
Computer Science Dept.



Uday Kamath, Ph.D.
Digital Reasoning Systems, Inc.



Barney Bishop, Ph.D.
George Mason University
Chemistry Dept.



Manpriya Dua, M.S.
George Mason University
Computer Science Dept.

*Special thanks to members of NIAID BCBB, the Shehu lab and Jianlin Cheng (U. Missouri) for their helpful feedback and suggestions.*

NIH > National Institute of Allergy and Infectious Diseases

*Thank you for listening!*

*Questions?*

National Institute of
Allergy and
Infectious Diseases

NIH

NIAID

# Extra Slides

# Antimicrobial Resistance Rates



**ANTIBIOTIC RESISTANCE INDENTIFIED** / **ANTIBIOTIC INTRODUCED**

- penicillin-R Staphylococcus — 1940
- 1943 penicillin
- 1950 tetracycline
- 1953 erythromycin
- tetracycline-R Shigella — 1959
- 1960 methicillin
- methicillin-R Staphylococcus — 1962
- penicillin-R pneumococcus — 1965
- 1967 gentamicin
- erythromycin-R Streptococcus — 1968
- 1972 vancomycin
- gentamicin-R Enterococcus — 1979
- 1985 imipenem and ceftazidime
- ceftazidime-R Enterobacteriaceae — 1987
- vancomycin-R Enterococcus — 1988
- levofloxacin-R pneumococcus — 1996
- 1996 levofloxacin
- imipenem-R Enterobacteriaceae — 1998
- XDR tuberculosis — 2000
- 2000 linezolid
- linezolid-R Staphylococcus — 2001
- vancomycin-R Staphylococcus — 2002
- 2003 daptomycin
- PDR-Acinetobacter and Pseudomonas — 2004/5
- ceftriaxone-R Neisseria gonorrhoeae PDR-Enterobacteriaceae — 2009
- 2010 ceftaroline
- ceftaroline-R Staphylococcus — 2011

Source: US Center for Disease Control
https://www.cdc.gov/drugresistance/about.html

National Institute of
Allergy and
Infectious Diseases

NIAID

28

# Learning Curves (10-Fold CV)

# Performance Comparison on Other AMP Data Sets

| Method | Data Set | No. AMPs (Overlap) | No. Non-AMPs (Overlap) | ACC(%) | MCC |
|--------|----------|--------------------|------------------------|--------|-----|
| Our DNN | | | | **92.95** | **0.860** |
| AntiBP2 | | | | 91.64 | 0.831 |
| CAMP ANN | Lata et al. 2010 | 999 (75%) | 999 (0%) | 81.03 | 0.624 |
| CAMP DA | | | | 84.28 | 0.690 |
| CAMP RF | | | | 87.09 | 0.752 |
| CAMP SVM | | | | 86.69 | 0.739 |
| iAMP-2L | | | | 86.34 | 0.735 |
| iAMPpred | | | | 92.84 | 0.858 |
| Our DNN | | | | **90.93** | **0.827** |
| AntiBP2 | | | | 85.30 | 0.706 |
| CAMP ANN | Fernandes et al. 2012 | 115 (62%) | 116 (0%) | 77.06 | 0.553 |
| CAMP DA | | | | 77.06 | 0.572 |
| CAMP RF | | | | 79.65 | 0.640 |
| CAMP SVM | | | | 77.06 | 0.584 |
| iAMP-2L | | | | 87.90 | 0.759 |
| iAMPpred | | | | 84.00 | 0.691 |
| Our DNN | | | | **97.42** | **0.949** |
| AntiBP2 | | | | 89.10 | 0.781 |
| CAMP ANN | Xiao et al. 2013 | Train Set: 878 (77%) | Train Set: 2368[†] (0.3%) | 80.00 | 0.610 |
| CAMP DA | | | | 71.79 | 0.487 |
| CAMP RF | | Test Set: 920 (62%) | Test Set: 920 (0%) | 65.27 | 0.396 |
| CAMP SVM | | | | 67.77 | 0.429 |
| iAMP-2L | | | | 92.23 | 0.845 |
| iAMPpred | | | | 72.99 | 0.509 |

†37 sequences were removed from the original data set to remove duplicates or peptides containing fragments identical to known AMPs as in Veltri (2015).

National Institute of Allergy and Infectious Diseases

# Weights and Layers



Inputs X

Layer 1 Layer 2 ... Layer *m*

weights $L_1$    weights $L_2$    ...    weights $L_m$

Predictions Y

$$\begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix} \qquad \begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{bmatrix}$$

Often weights are randomly initialized and layer outputs are often "activated" using functions to force numbers in a certain range. Typical examples include: *sigmoid*, *tanh*, and **rectifier linear unit** (ReLU) functions.

National Institute of Allergy and Infectious Diseases

**W**



$$\begin{bmatrix} \vdots \\ x_n \end{bmatrix}$$

Inputs X

weights
$L_1$





$$\begin{bmatrix} \vdots \\ y_n \end{bmatrix}$$

Predictions Y

Often w...             ...s are often
"activated...          ...ange. Typical
exampl...              ...**unit** (ReLU)

NIH
National Institute of
Allergy and
Infectious Diseases

NIAID

# Weights and Layers – Optimizing the Network



Start Round 2 of Training!

Layer 1   Layer 2   . . .   Layer $m$

$$\begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{bmatrix}$$

Inputs **X**

weights $L_1$   weights $L_2$   . . . weights $L_m$

$$\begin{bmatrix} y_1 \\ y_2 \\ \dots \\ y_n \end{bmatrix}$$

Predictions **Y'**

Optimizer (adjust weights)

Apply scoring function using *actual* **Y** values

National Institute of Allergy and Infectious Diseases

# Optimizing via Backpropagation
## *The Secret Sauce of Deep Neural Networks (DNNs)*

- How do DNNs learn so well? The key is they compute answers across layers in a *forward pass* and then to use a *backwards pass* to optimize the weights. This way **all** layers are updated each round (sometimes called an 'epoch') of training!

- How does this backward pass work? **The chain rule** (remember from calculus?) – where we can calculate the derivative (the slope or rate of change) from *two or more* functions.

You might have seen this written as: $(f \circ g)' = (f' \circ g) \cdot g'$

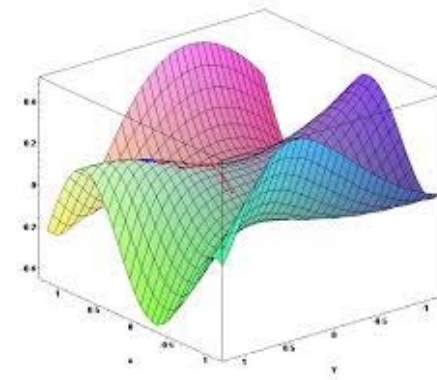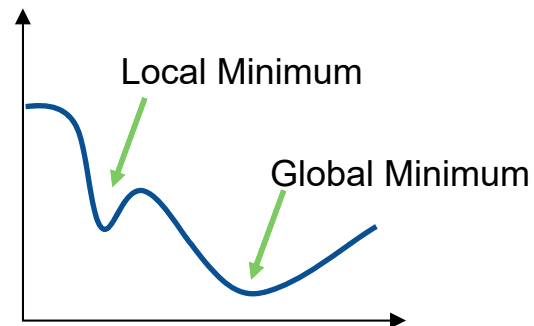or maybe like this: $\dfrac{dz}{dx} = \dfrac{dz}{dy} \cdot \dfrac{dy}{dx}$   or maybe like this: $F'(x) = f'\big(g(x)\big)g'(x)$

*The takeaway: we can calculate the derivative using multiple functions at the same time!*

National Institute of Allergy and Infectious Diseases

# Optimizing via Backpropagation
## *The Secret Sauce of Deep Neural Networks (DNNs)*



- For our DNNs we are calculating the *gradient* (a vector of derivatives) to account for the change across the network based on the forward pass results.

- Given a function $f(x)$ where x's are our **training** inputs- the gradient forms a vector: $\nabla f(x) =$

$$[ \frac{\partial f}{\partial x} , \frac{\partial f}{\partial y} ] = [y, x]$$



Local Minimum

Global Minimum



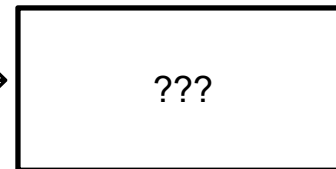National Institute of Allergy and Infectious Diseases

NIAID

# Some Backpropagation Intuition

Lets look at multiplication:

$$f(x,y) = xy \quad \rightarrow \frac{\partial f}{\partial x} = y \quad \frac{\partial f}{\partial y} = x$$

**If $x$ = 4 and $y$ = −3**

$$f(x,y) = -12 \qquad \frac{\partial f}{\partial x} = -3 \quad \frac{\partial f}{\partial y} = 4$$

Lets look at basic addition: $\quad f(x,y) = x + y \quad \rightarrow \boxed{\text{???}}$

*What happens to each function if we change $x$*
*… or change $y$?*

National Institute of
Allergy and
Infectious Diseases

# Some Backpropagation Intuition
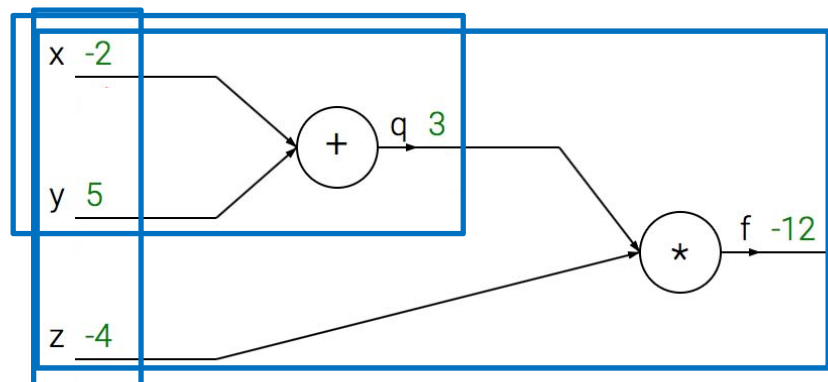
Lets look at **multiple** functions:

$$f(x, y, z) = (x + y)\,z$$

We can rewrite this as: $q = x + y \ \text{ and } f = qz$

SO $\dfrac{\partial f}{\partial q}$ = z , $\dfrac{\partial f}{\partial z}$ = q    … for $(x + y)$ as we saw before: $\dfrac{\partial f}{\partial x}$ = 1 , $\dfrac{\partial f}{\partial y}$ = 1

The **chain rule** says **multiply**: $\dfrac{\partial f}{\partial x} = \dfrac{\partial f}{\partial q} \cdot \dfrac{\partial q}{\partial x}$

National Institute of
Allergy and
Infectious Diseases

*Lets look at this with code and a visual representation!*

# Backpropagation Example



```
# set some inputs
x = -2; y = 5; z = -4

# perform the forward pass
q = x + y # q becomes 3
f = q * z # f becomes -12

# perform the backward pass (backpropagation)
# first backprop through f = q * z
dfdq = z # df/dq = z, so gradient on q becomes -4
dfdz = q # df/dz = q, so gradient on z becomes 3

# now backprop through q = x + y
dfdx = 1.0 * dfdq # dq/dx = 1  chain rule!
dfdy = 1.0 * dfdq # dq/dy = 1
```

x  -2

y  5

q  3

z  -4

f  -12

+

*

*Forward pass is* **green**. *Backward pass is* **red**.

**The big takeaway:** The final gradient [ $\frac{\partial f}{\partial x}$ , $\frac{\partial f}{\partial y}$ , $\frac{\partial f}{\partial z}$ ] tells us how sensitive

our function $f$ is to the variables $x$ , $y$ , and $z$.

National Institute of
Allergy and
Infectious Diseases