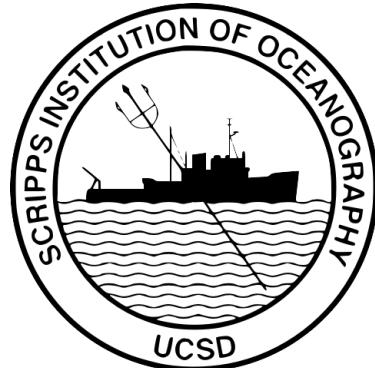




Gary's Unbelievable Research Unit



Using Deep Siamese Neural Networks to Speed Up Natural Products Research

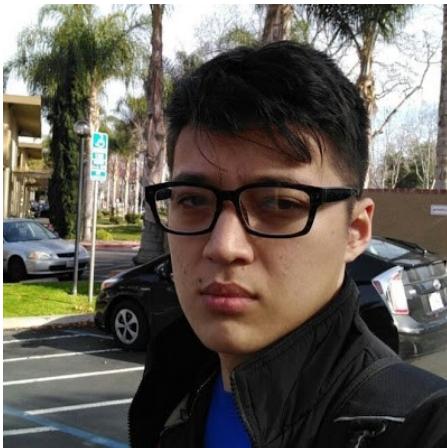
Garrison W. Cottrell

CSE Department

Gary's Unbelievable Research Unit (GURU)

UC San Diego

Collaborators



Yerlan Idelbayev
UC Merced



Yash_Nannapaneni
UCSD

1/19/19

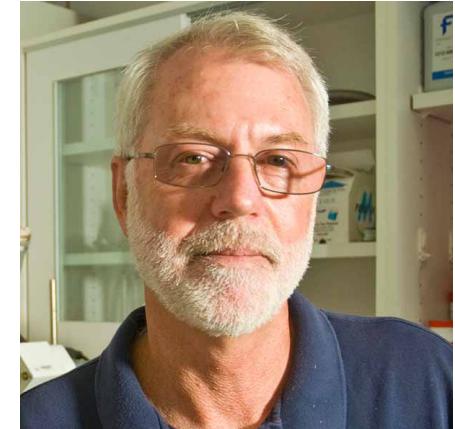


Chen Zhang
Scripps Institution
of Oceanography



Nick Roberts, UCSD

Southern China University of Technology



Bill Gerwick
Scripps Institution
of Oceanography



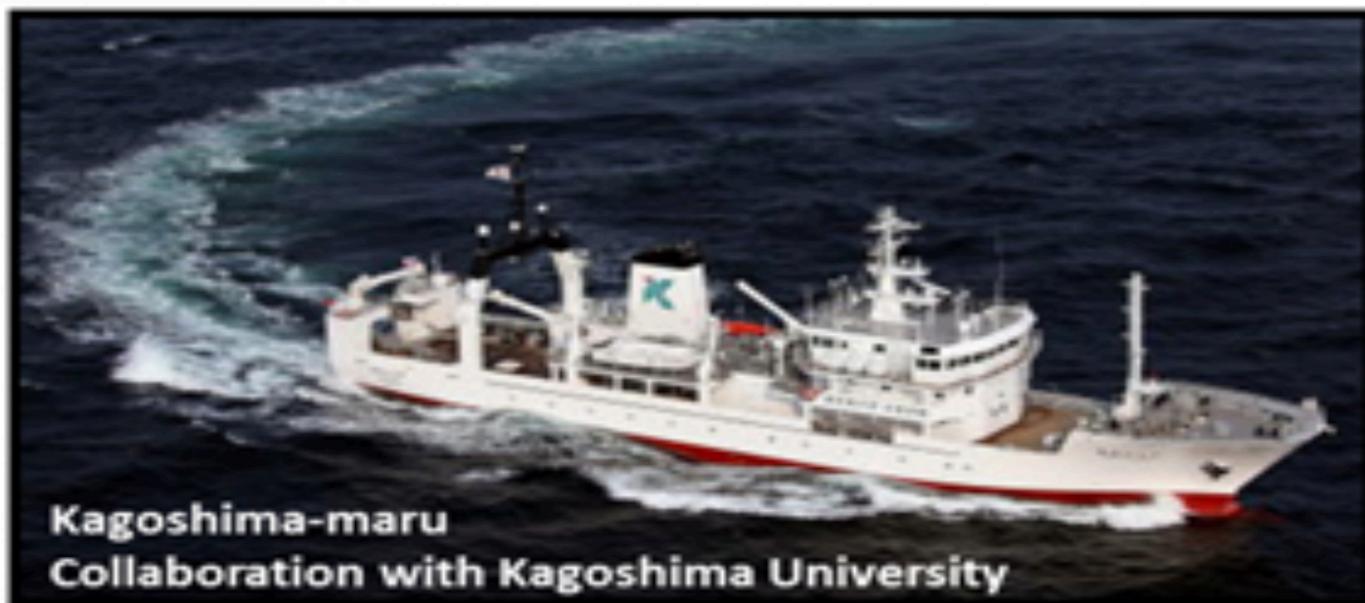
**Poornav Sargur
Purushothama**
UCSD

Outline

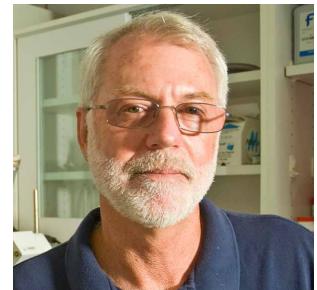
- Natural products research
 - What is it
 - Why is it important
- Deep Siamese Convolutional Networks
 - What are they
 - Why are they cool
- SMART
 - Methods
 - Results

Natural Products Research

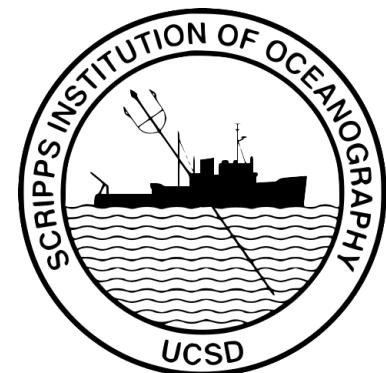
- 50% the approved drugs during the last 30 years are derived from natural products (Newman and Cragg 2012)



Natural Products Research



Bill Gerwick
Scripps Institution
of Oceanography



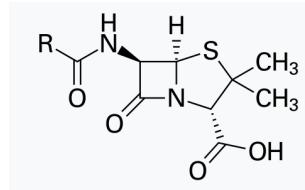
Here's Bill himself collecting samples in the Palmyra Atoll



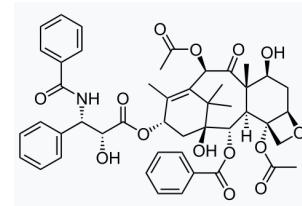
Natural Products Research

- 50% the approved drugs during the last 30 years are derived from natural products (Newman and Cragg 2012)

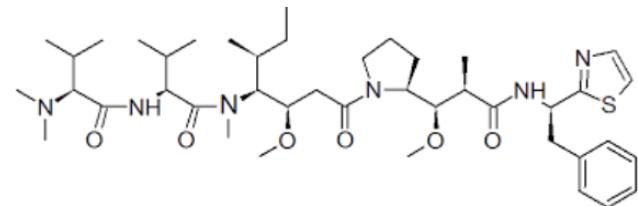
- Penicillin (1928)



- Taxol (1971) an anti-cancer agent derived from Yew tree bark



- Dolastatin 10 (2001) another anti-cancer agent derived from a sea hare that feeds on tropical cyanobacteria in the Indian Ocean, and later traced back to the cyanobacterium itself.



Natural Products Research is ongoing -and a lot of it is about structure discovery

Article

Five macrocyclic glycosides from *Schoenoplectus tabernaemontani* >

Dian Peng, Xiaolin Li

Published online: 30

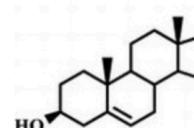


Article

Synthesis, characterization, and biological evaluations
of some steryl 2-methoxybenzoates as anticancer
agents >

Yanmin Huang, Hai
Xiaolan Liu, Junyan

Published online: 30



Article
A new flavonoid from the leaves of *Atalantia monophylla* (L.) DC >

Priyapan Posri, Jittra Suthiwong
Chindawadee Chuenban, Chant

Published online: 30 Mar 2018



Article

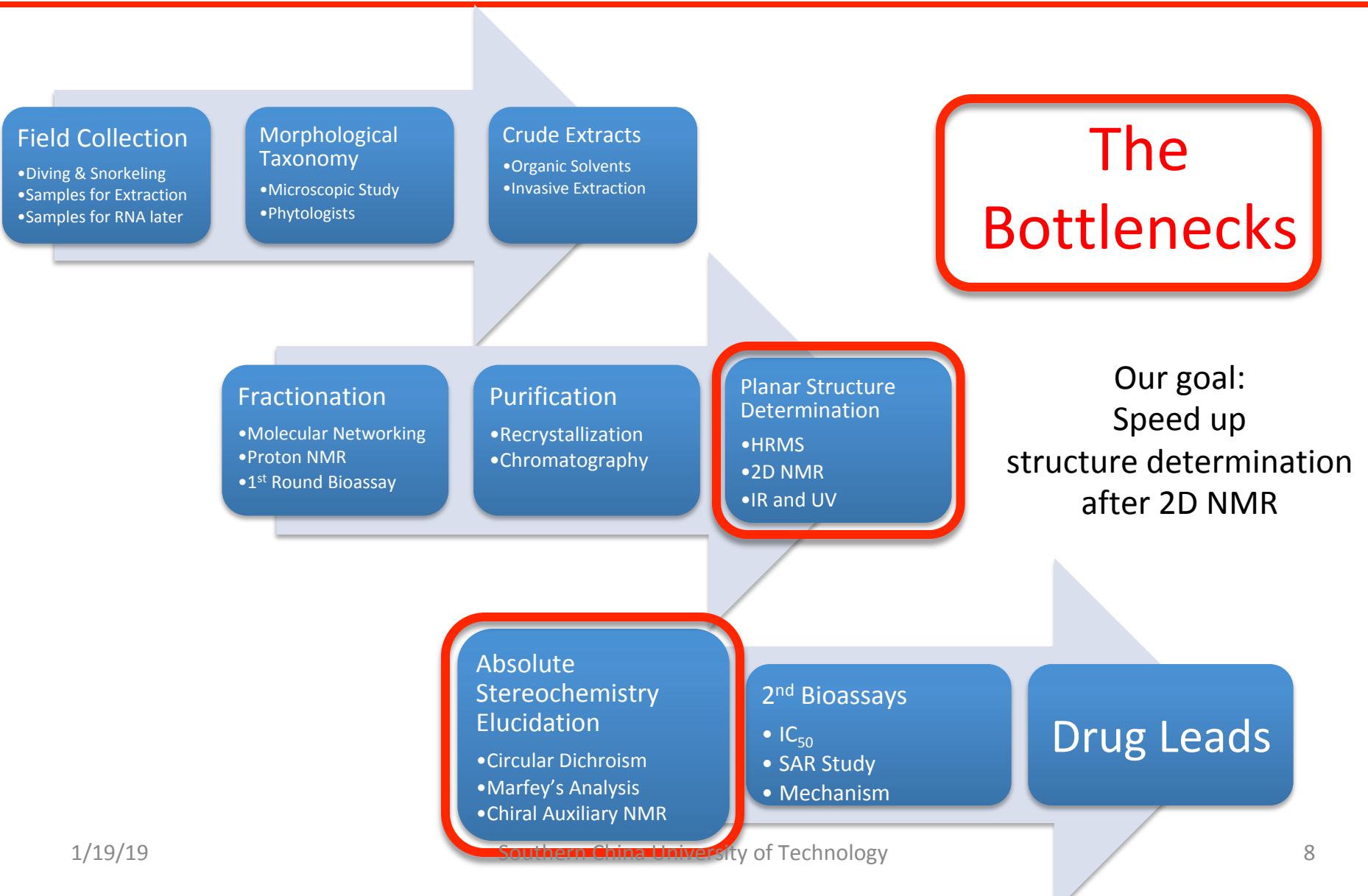
Synthesis and antiproliferative activities of
thioxoflavonoids on three human cancer cells >

Wei Li, Peipei Han, Shuanglian Cai & Qiuwan Wang

Published online: 27 Mar 2018



Current Natural Products Discovery is a Slow Process



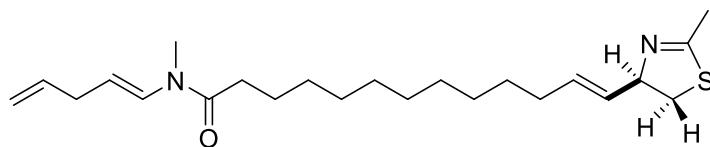
Natural Products Research

- Figuring out the *structure* of any new compound starts with the 2D NMR spectrum of the compound – the “fingerprint” of the compound
- Using a deep, convolutional siamese neural network, we have developed a system called SMART that *clusters* compounds directly from their NMR spectrum.
- This allows new compounds to be mapped into the cluster space, giving an idea of what compounds it is similar to.
- This can speed up structure discovery

2D Nuclear Magnetic Resonance Imaging (2D NMR)

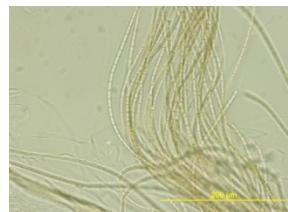
- 2D NMR is a technique that gives indications of bonds between atoms in a molecule
- We think of it as the “fingerprint” of the molecule
- And now I’ve told you everything I know about NMR!!

2D NMR is an Indicator of Compound Structure

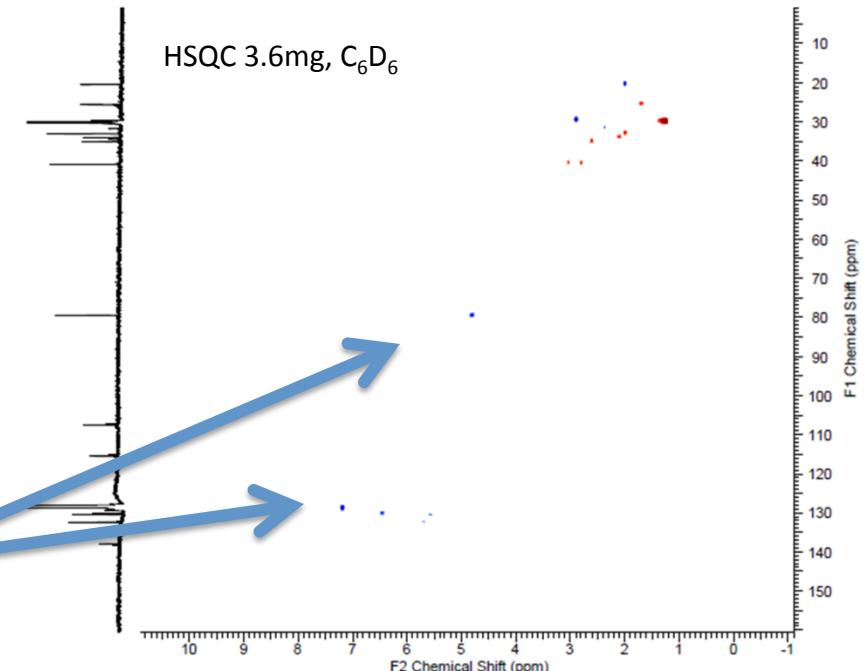
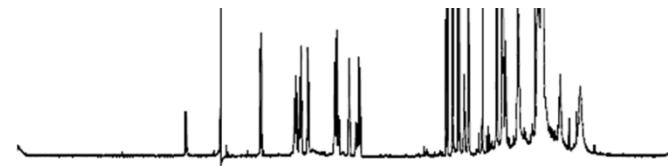


lauysteinamide A

Caldora penicillata

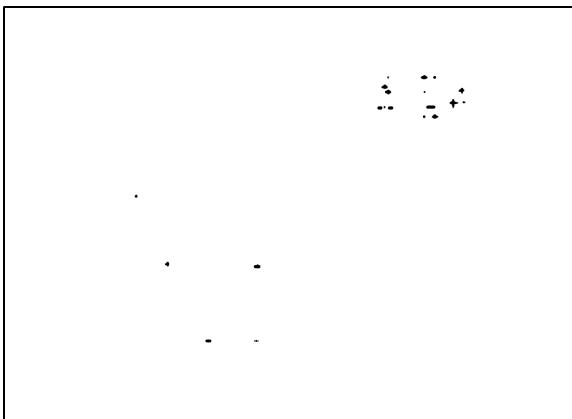


Each “dot” here corresponds to a bond

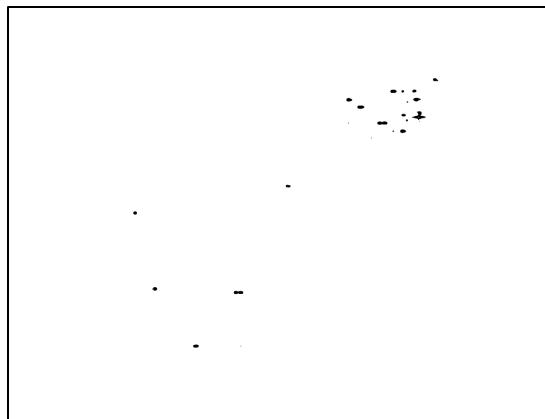


NMR structure analysis is subtle

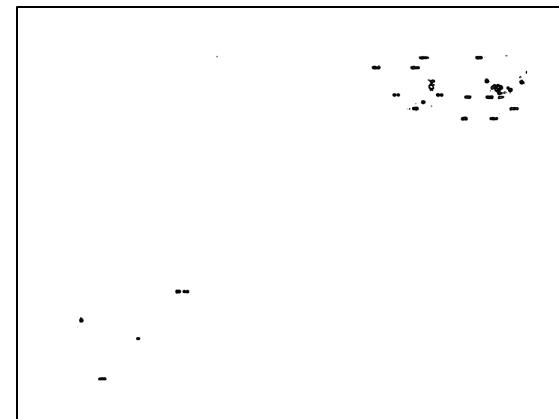
Aspewentin A



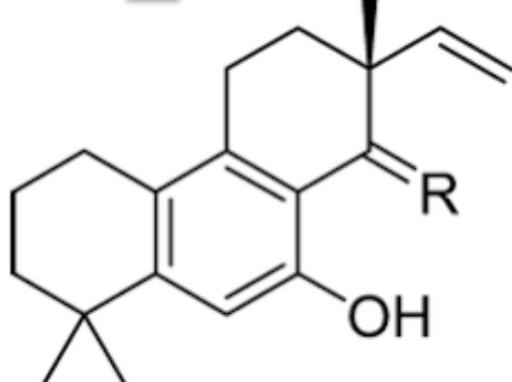
Aspewentin B



Aspewentin C



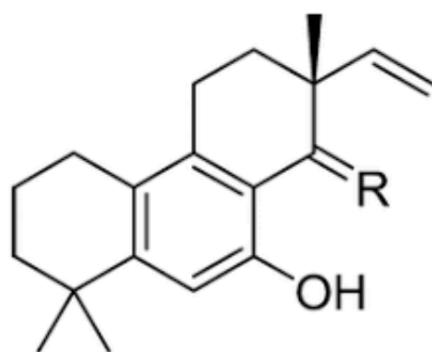
NMR Spectrum



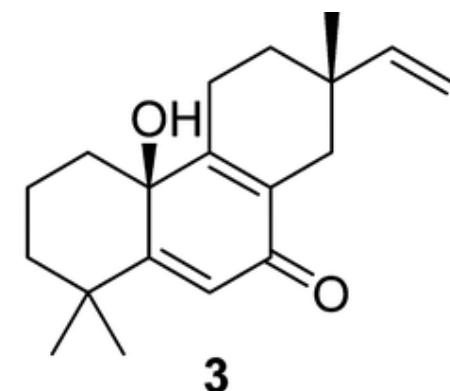
1 R = H₂



Actual Structure



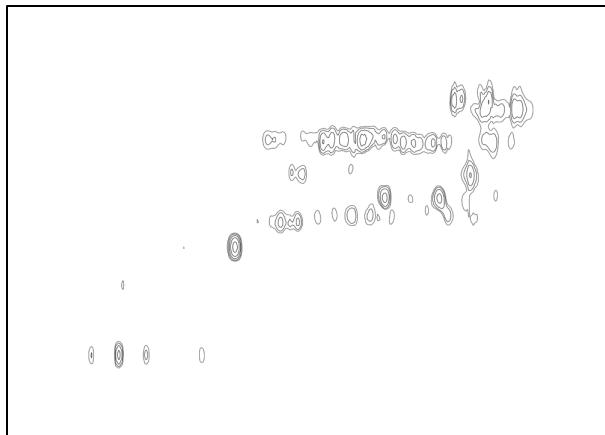
2 R= O



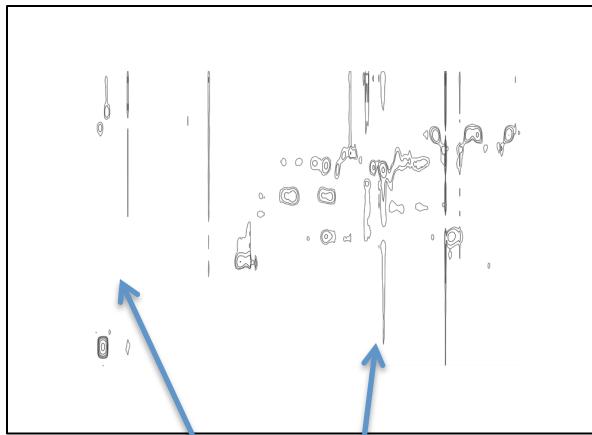
3

And Noisy!

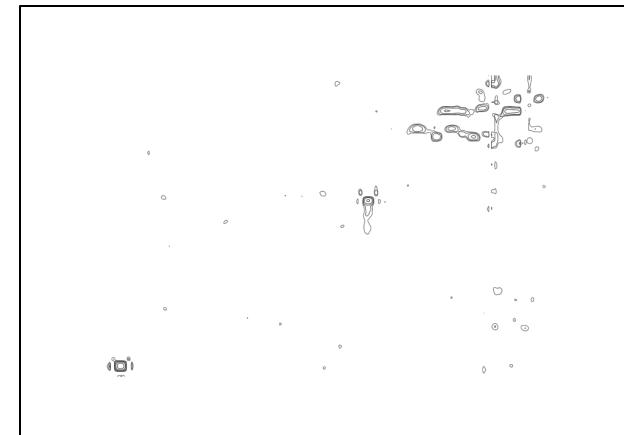
Chandonanone A



Chandonanone B



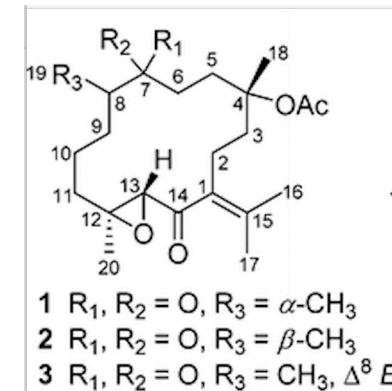
Chandonanone C



NMR Spectrum

Actual Structure

These are artifacts of the solvent used...

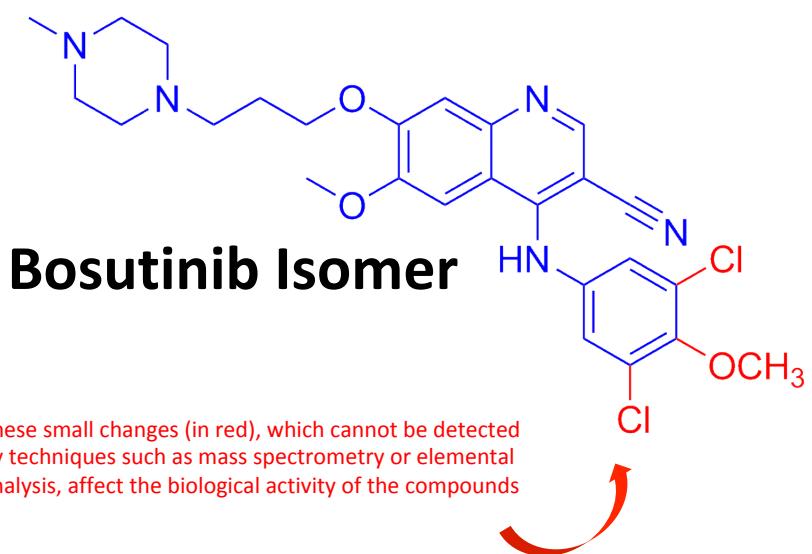
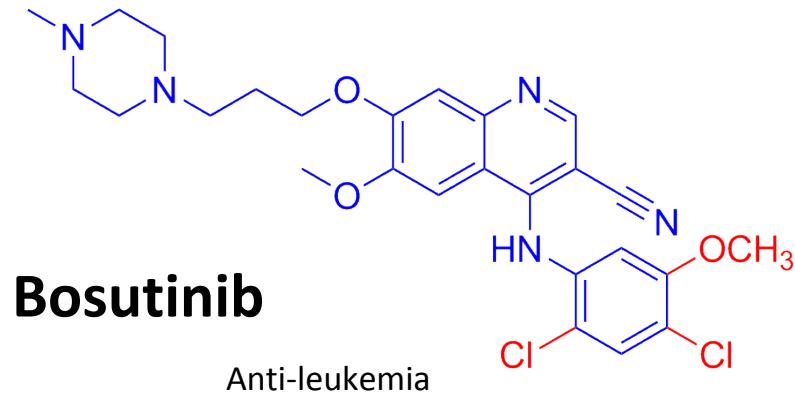


1 $R_1, R_2 = O, R_3 = \alpha\text{-CH}_3$

2 $R_1, R_2 = O, R_3 = \beta\text{-CH}_3$

3 $R_1, R_2 = O, R_3 = CH_3, \Delta^8 E$

And Mistakes can be Costly...



- Same Mass
- Same Chemical Composition
- Similar Chemical Structures
- NMR Misinterpreted by Human
- **Very Different Biological Function!**

"I still wonder what we could have done better. Clearly, not every lab can run and analyze an NMR spectrum of every chemical compound ordered."

— Dr. Ulrich Schweizer,
Researcher at Charité-Universitätsmedizin Berlin,
Germany

"We had wasted a huge amount of time and money on the wrong isomer."

— Dr. Steven G. Boxer,
Distinguished Professor of Chemistry at Stanford
University

Enter Deep Learning

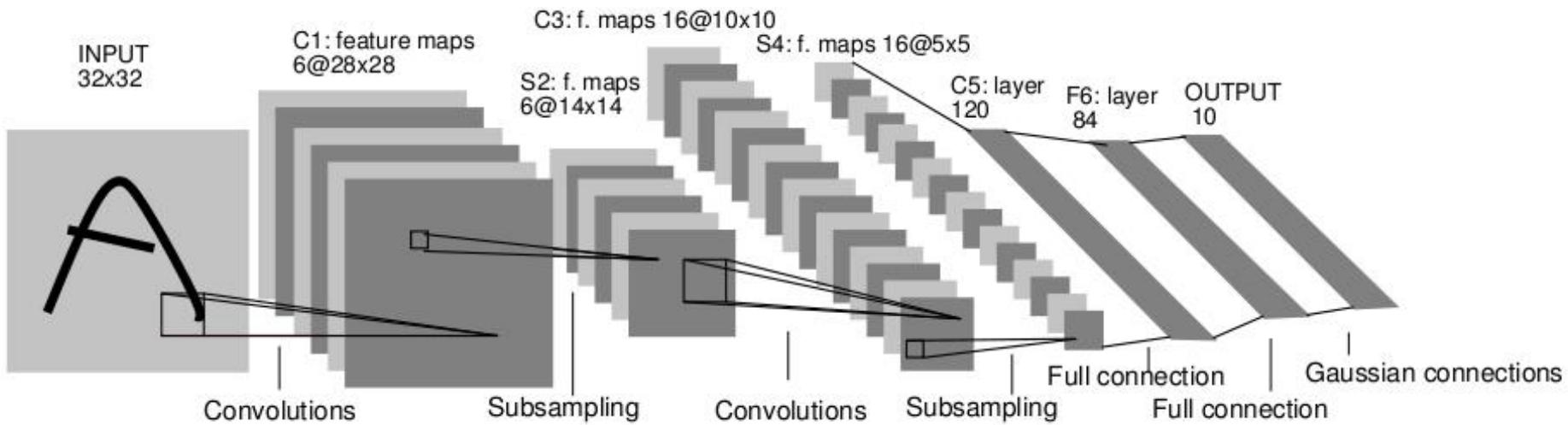
- Deep learning:
 - Any neural network with more than one hidden layer – usually many more!
- Convolutional Networks:
 - Have revolutionized computer vision since 2012
 - Learn many features
 - But apply the *same* feature across the image

Convolutional Neural Networks

Began with LeCun et al. 1989



1. Small, local receptive fields and learned features (kernels): *locality*
2. These are *replicated* across the image: *stationary statistics*
3. Spatial pooling: *translation invariance*



Convolutional Networks -- What they're good at: The Imagenet Large Scale Visual Recognition Challenge (ILSVRC)

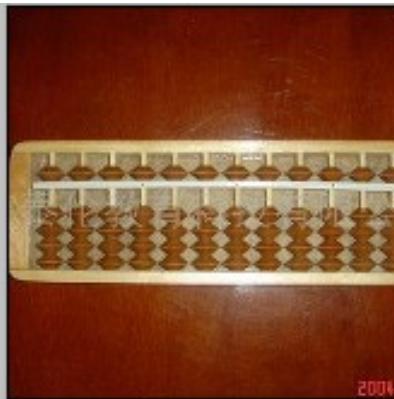
- 1.2 Million training images
- 1000 categories (732-1300 training images per class)
- 50,000 test images
- Large variation in images
- Many fine-scale categories (120 dog breeds, not just “dog”)

Imagenet challenge goal: Classification



lens cap

reflex camera
Polaroid camera
pencil sharpener
switch
combination lock



abacus

abacus
typewriter keyboard
space bar
computer keyboard
accordion



slug

slug
zucchini
ground beetle
common newt
water snake

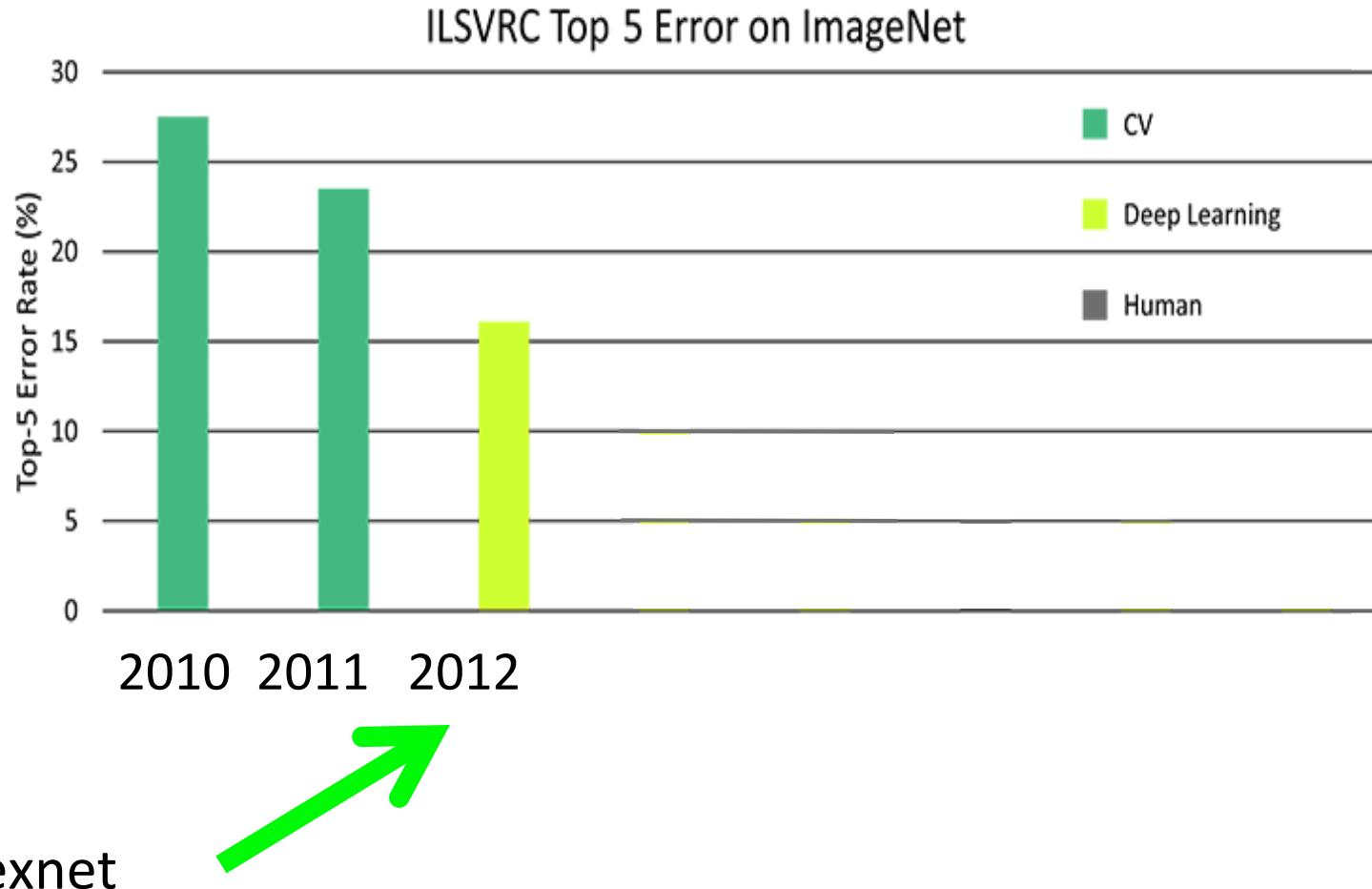


hen

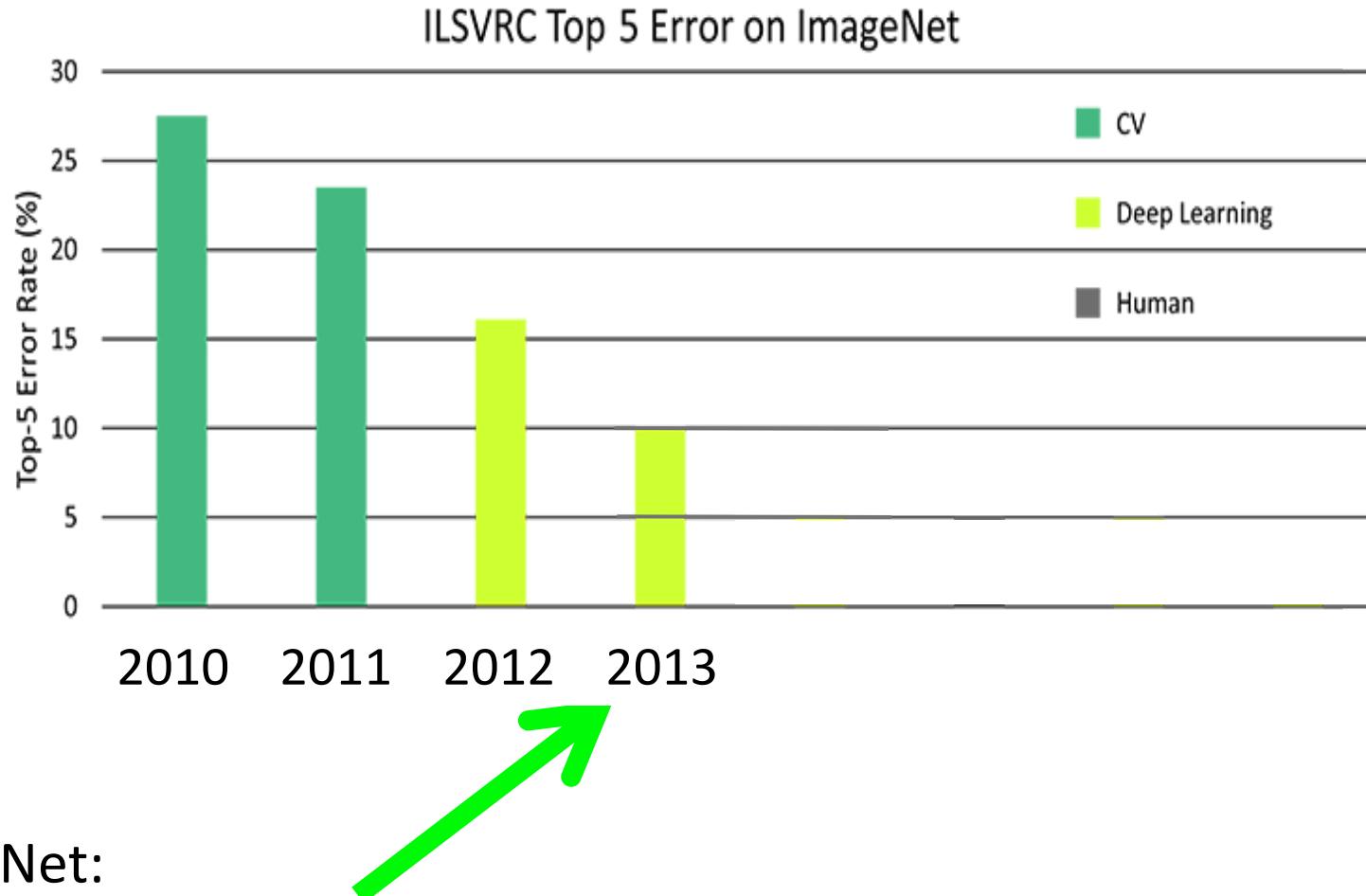
hen
cock
cocker spaniel
partridge
English setter

[Krizhevsky et al. NIPS 2012]

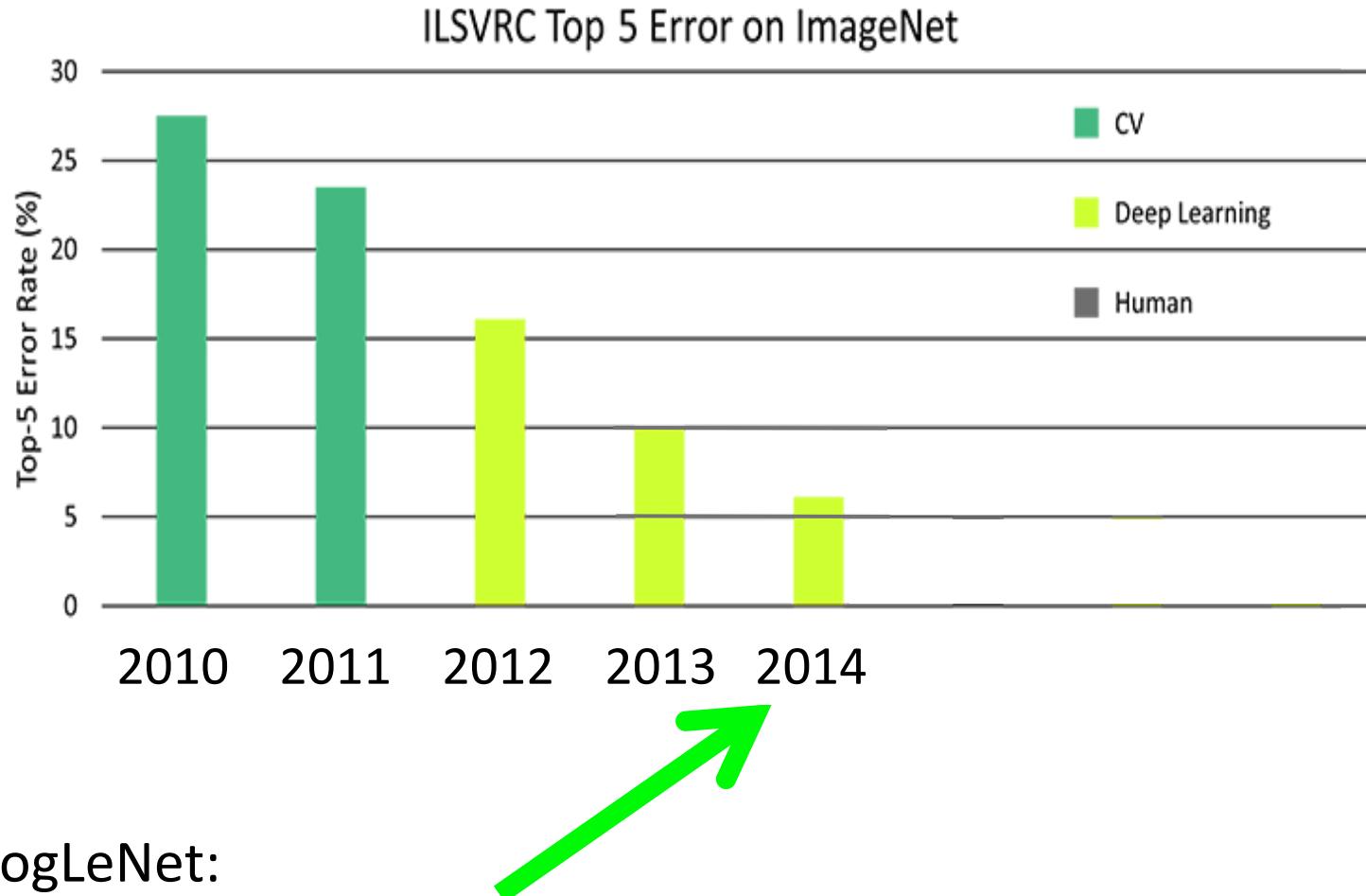
Why Deep Learning



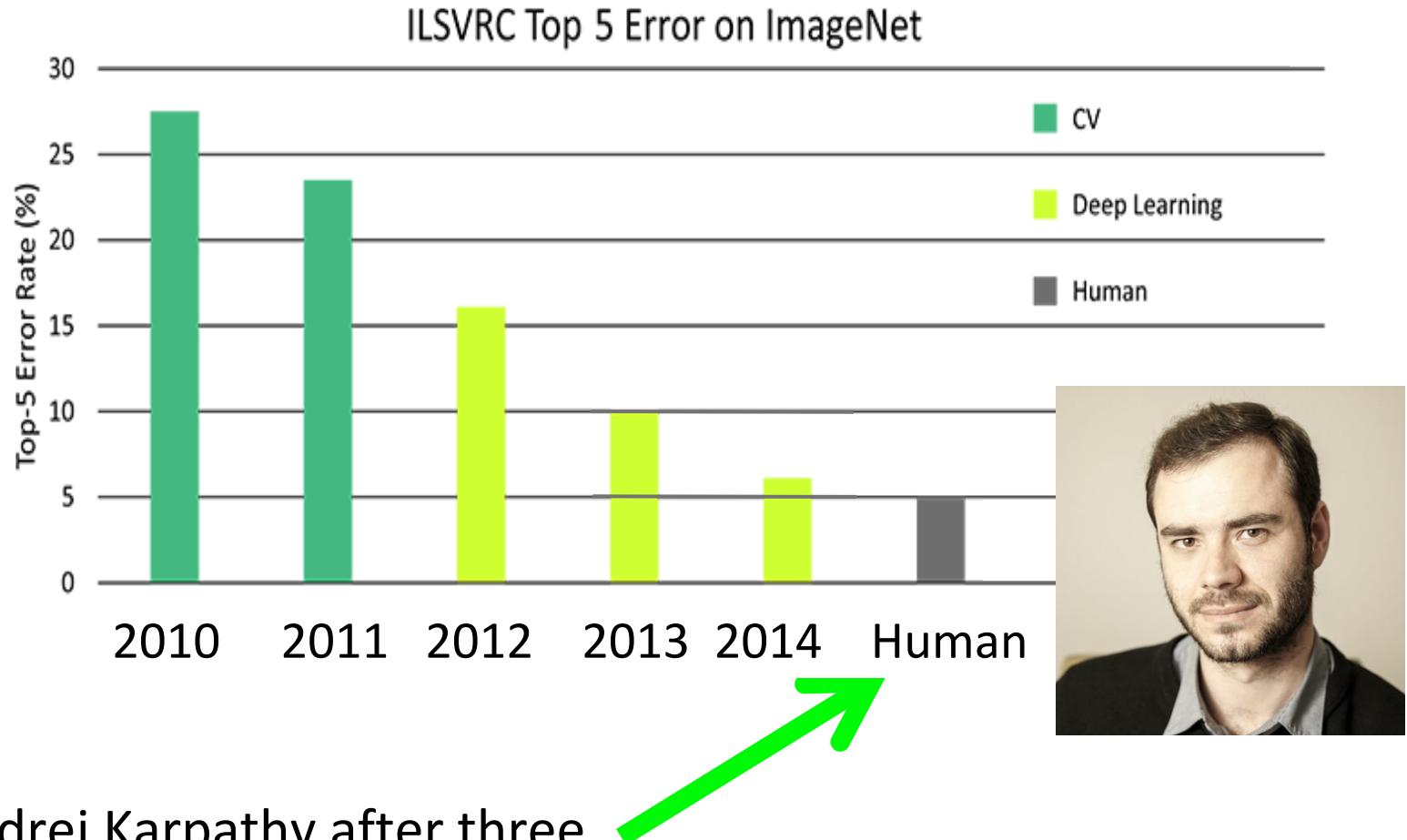
Why Deep Learning



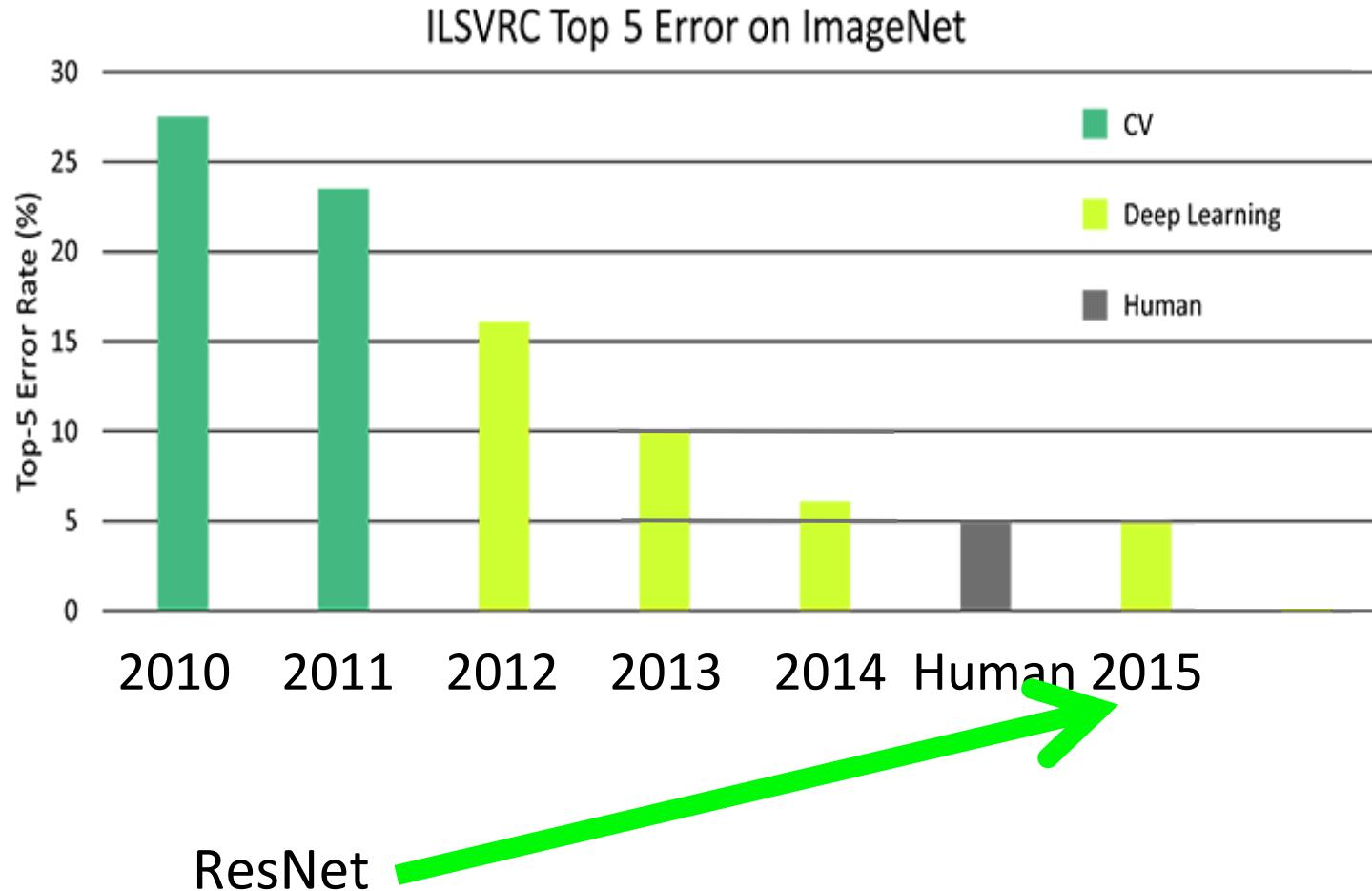
Why Deep Learning



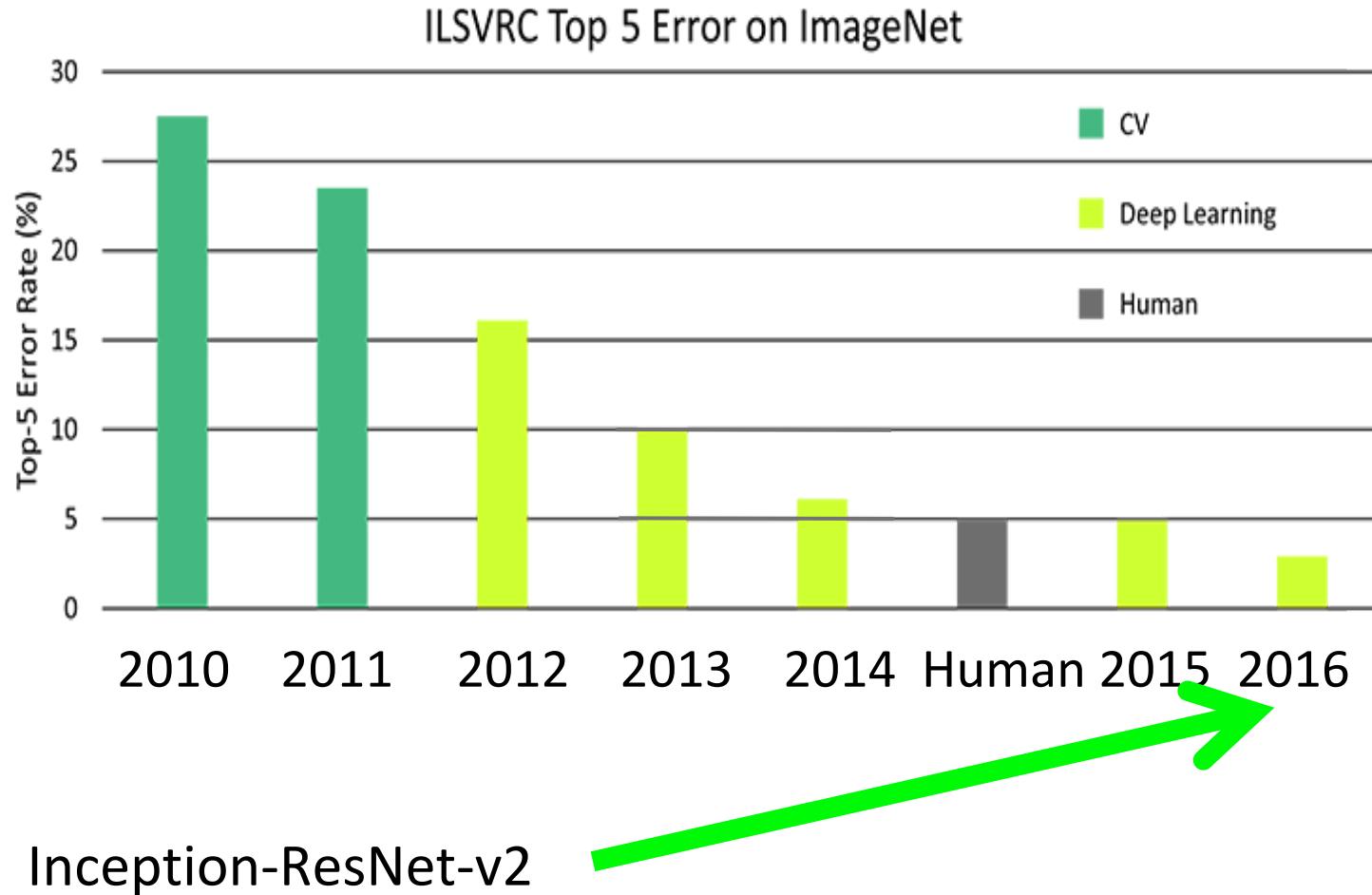
Why Deep Learning



Why Deep Learning

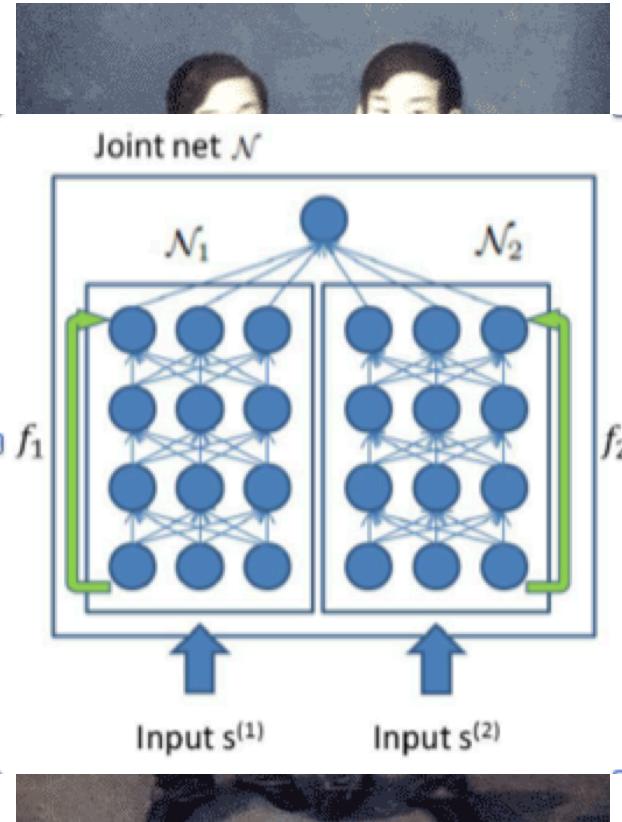


Why Deep Learning



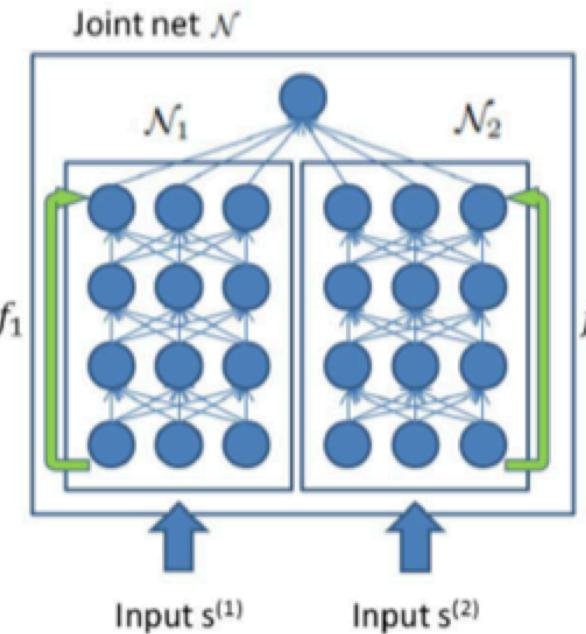
Siamese Neural Networks

- “Siamese twins” is a dated term for Conjoined Twins - twins that are physically joined at birth, sometimes sharing organs
- SNNs are the Neural Network analog for Conjoined Twins
- SNNs have two identical neural networks (like twins that are physically joined at birth)
- SNNs share their weights, both networks have the same weights (like twins sharing organs)
- They take two inputs and share a symmetric loss function (more about the loss function later)



Siamese Neural Networks

- Siamese Networks can be convolutional..
- They have linear output units, but they are not trained to any particular output – instead, they are trained to move outputs closer or farther apart.
- If two inputs are from the same category, minimize the distance between them
- If two inputs are from different categories, maximize the distance between them – up to some margin

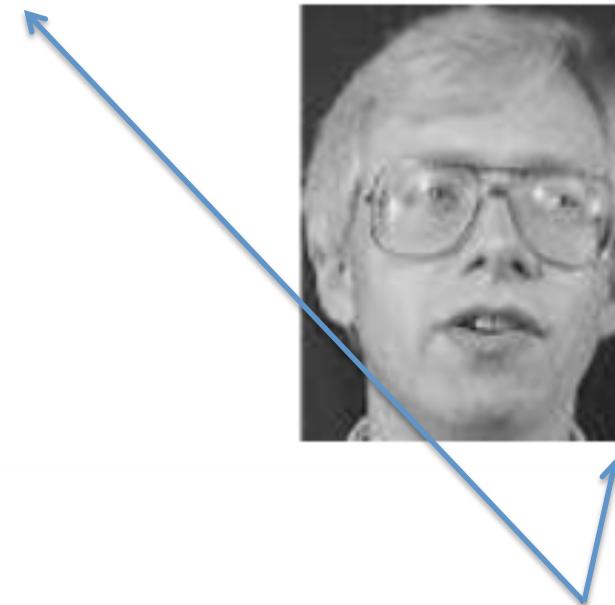


Same/Different Categories



Different:

Move the output vectors apart



Same:

Move the output vectors closer

Same/Different Categories

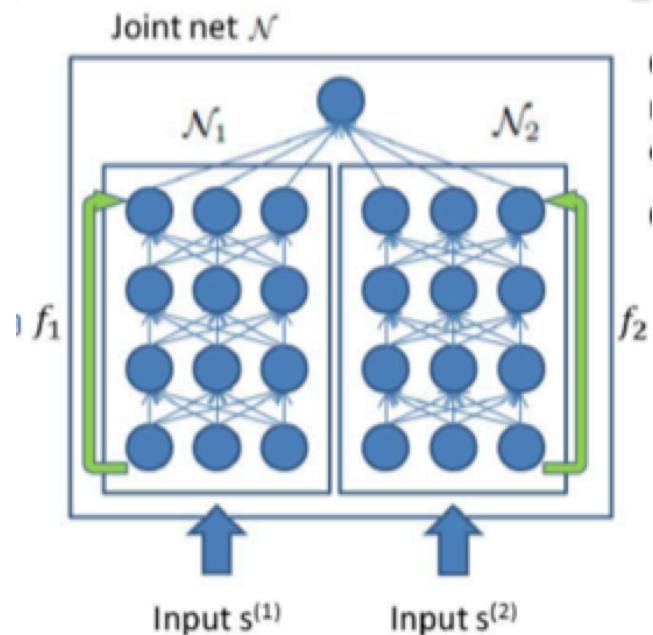


Note that putting all of these into the same category means the network has to learn to be *invariant* to the differences in these images

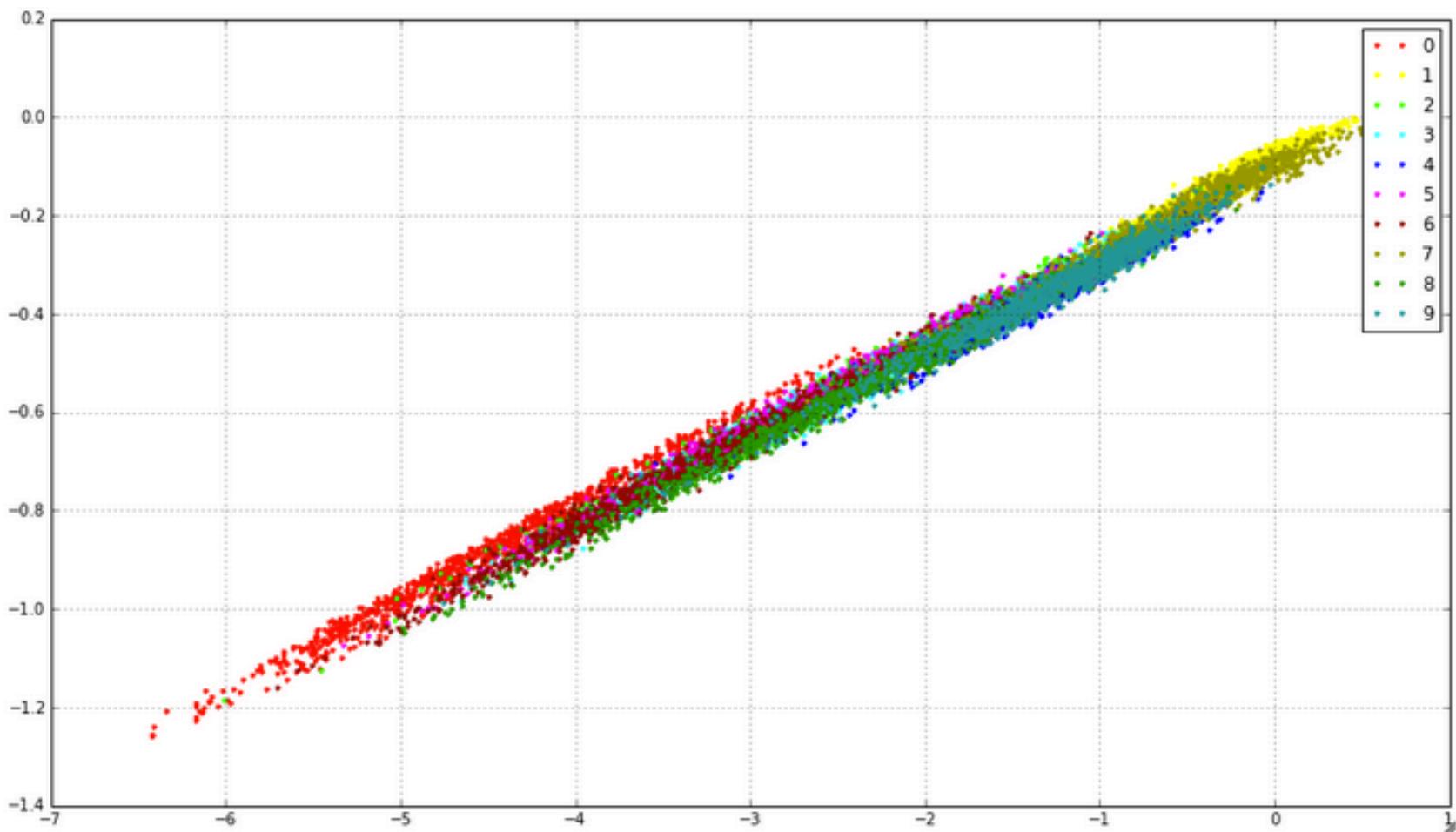
This transformation should *generalize* to new faces

Siamese Neural Networks

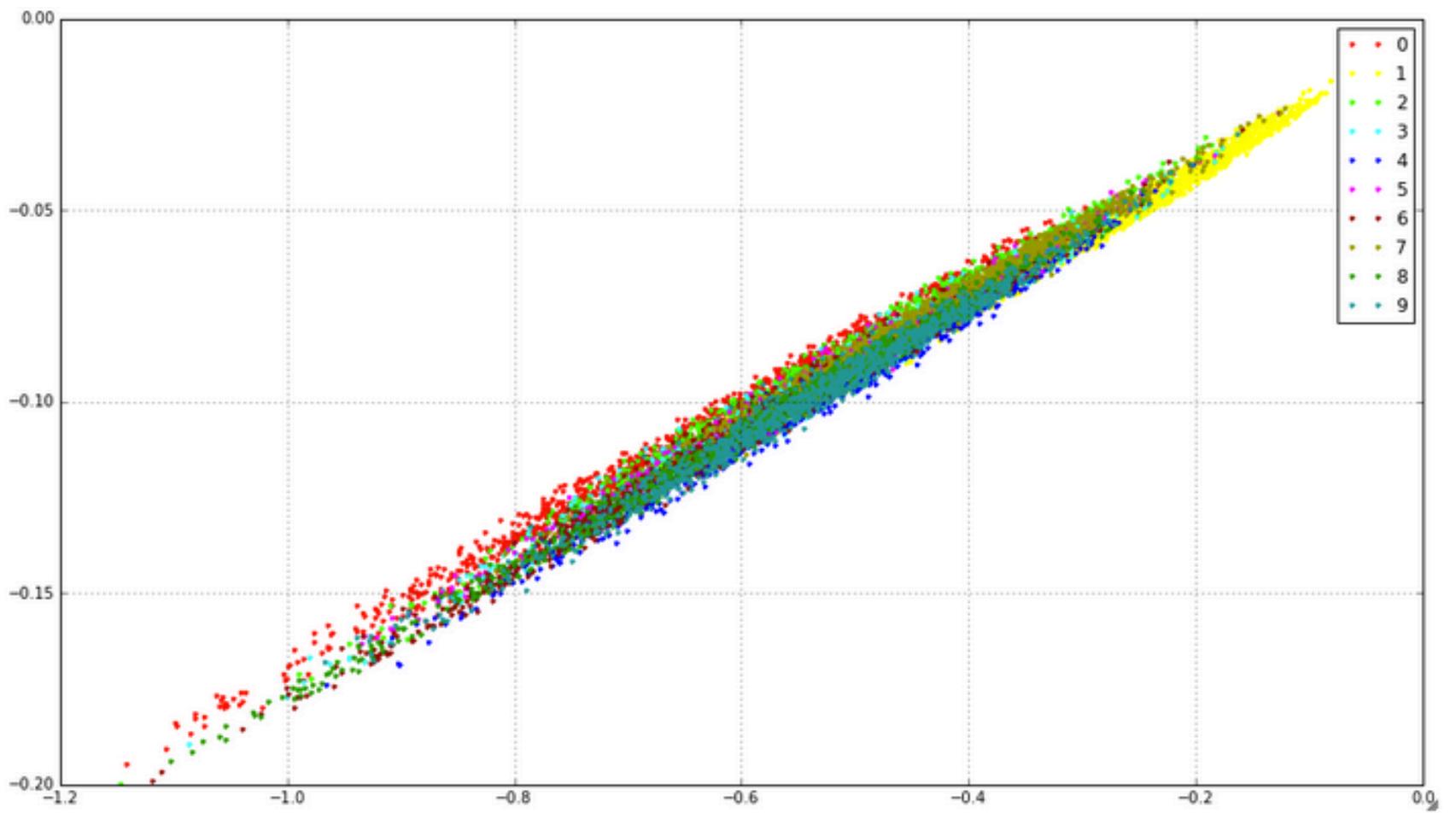
Next: MNIST Example using
two output units (so you can see
what they do



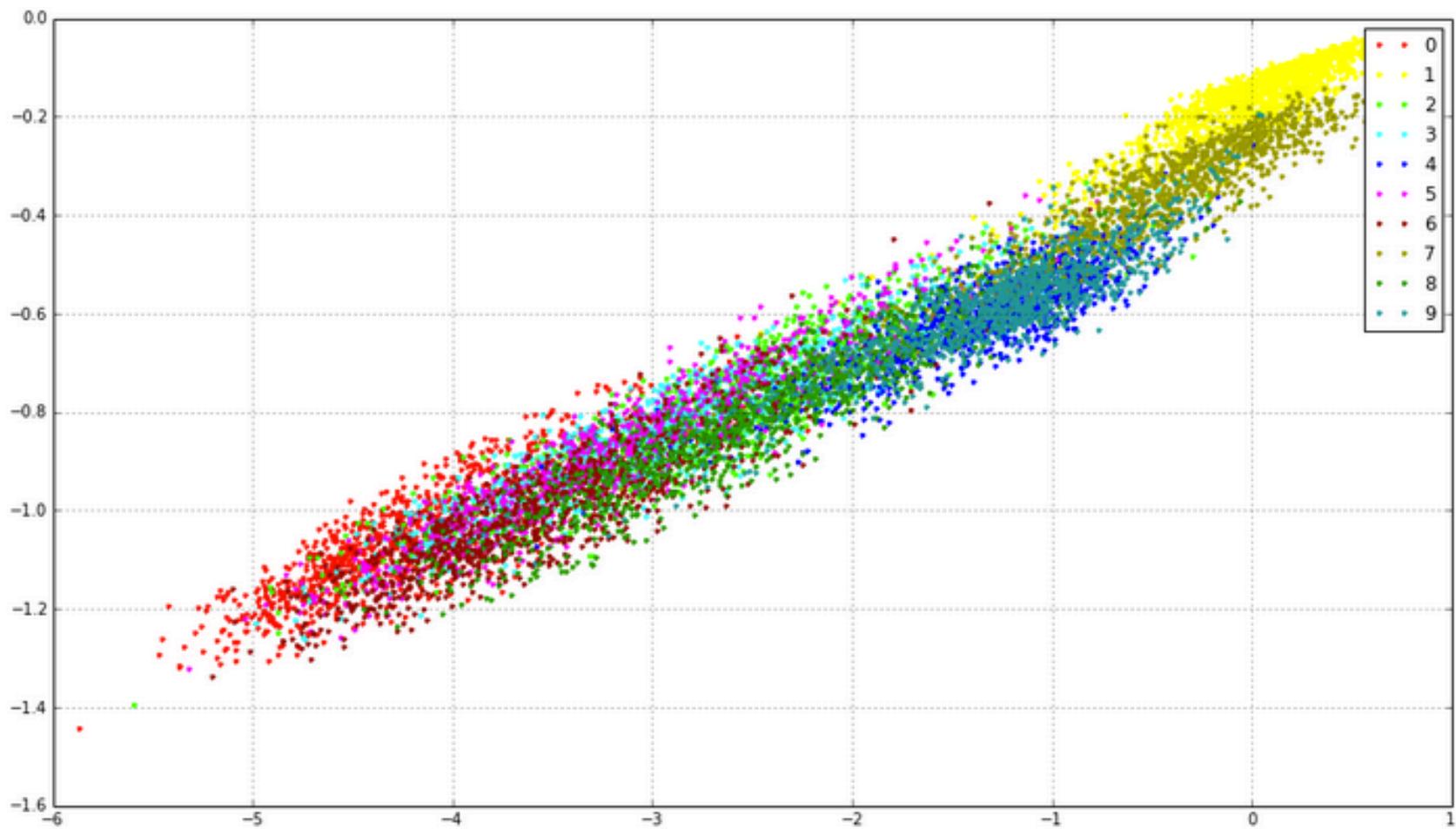
MNIST 100 iterations



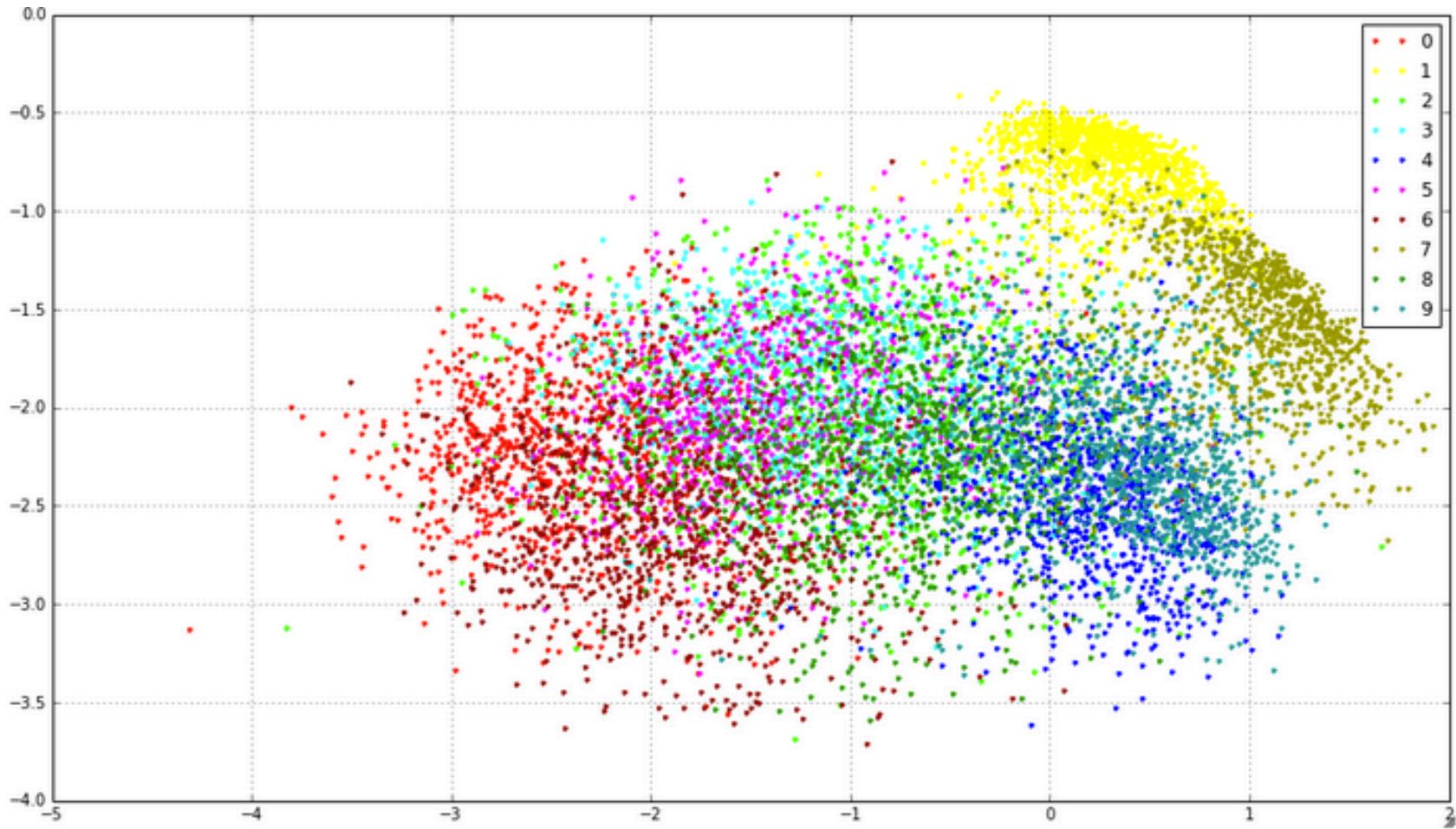
MNIST 200 iterations



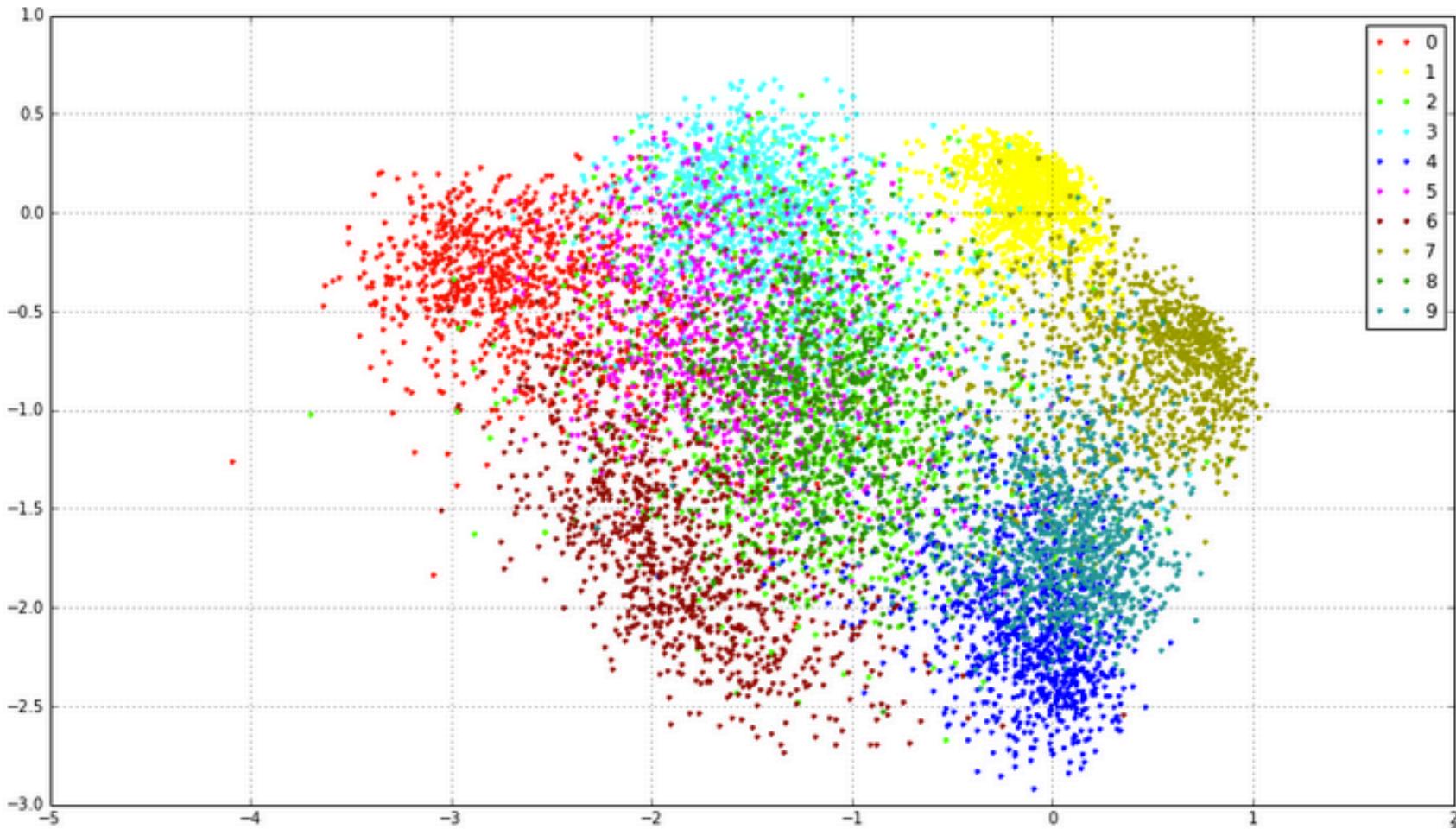
MNIST 300 iterations



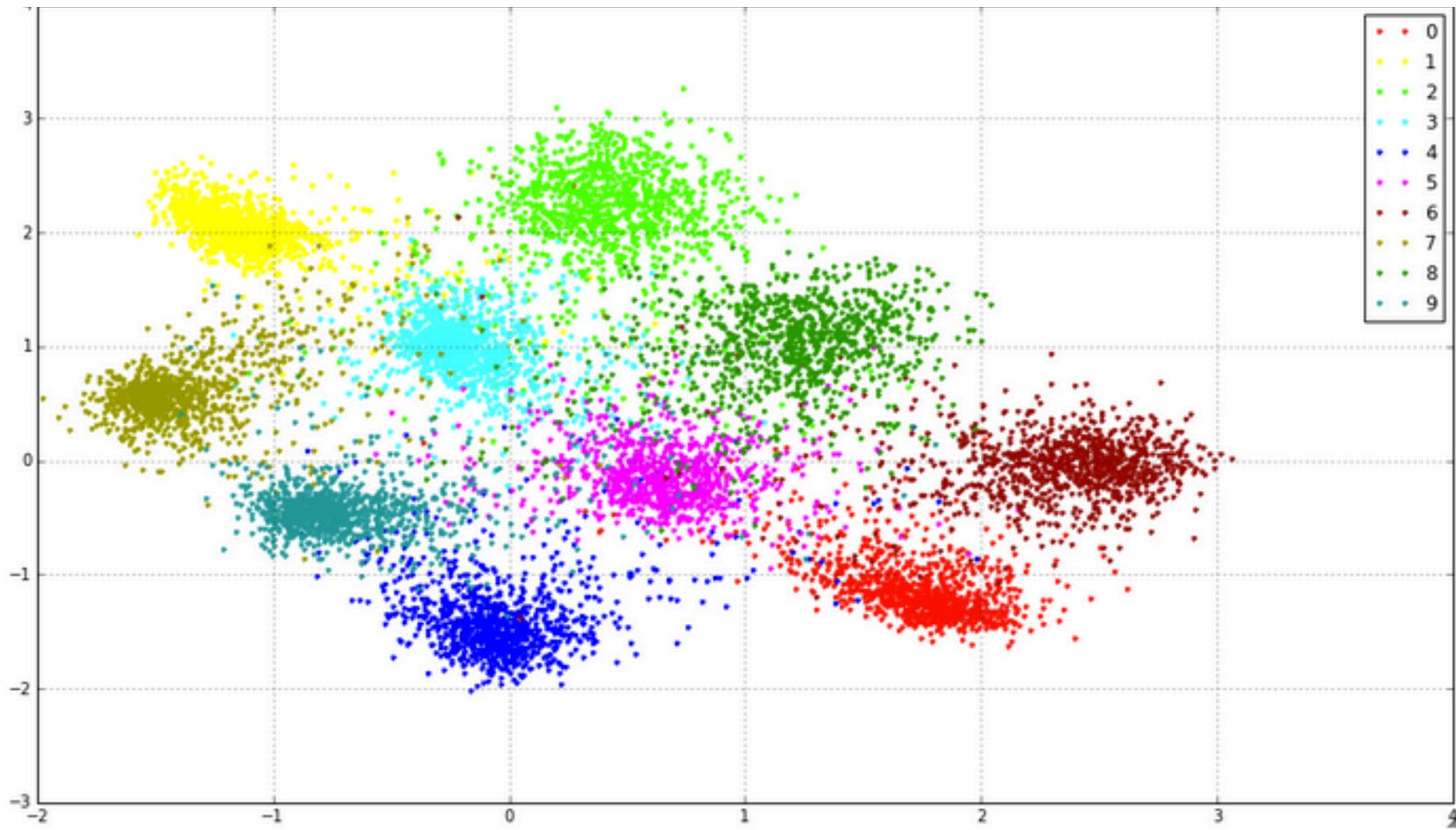
MNIST 400 iterations



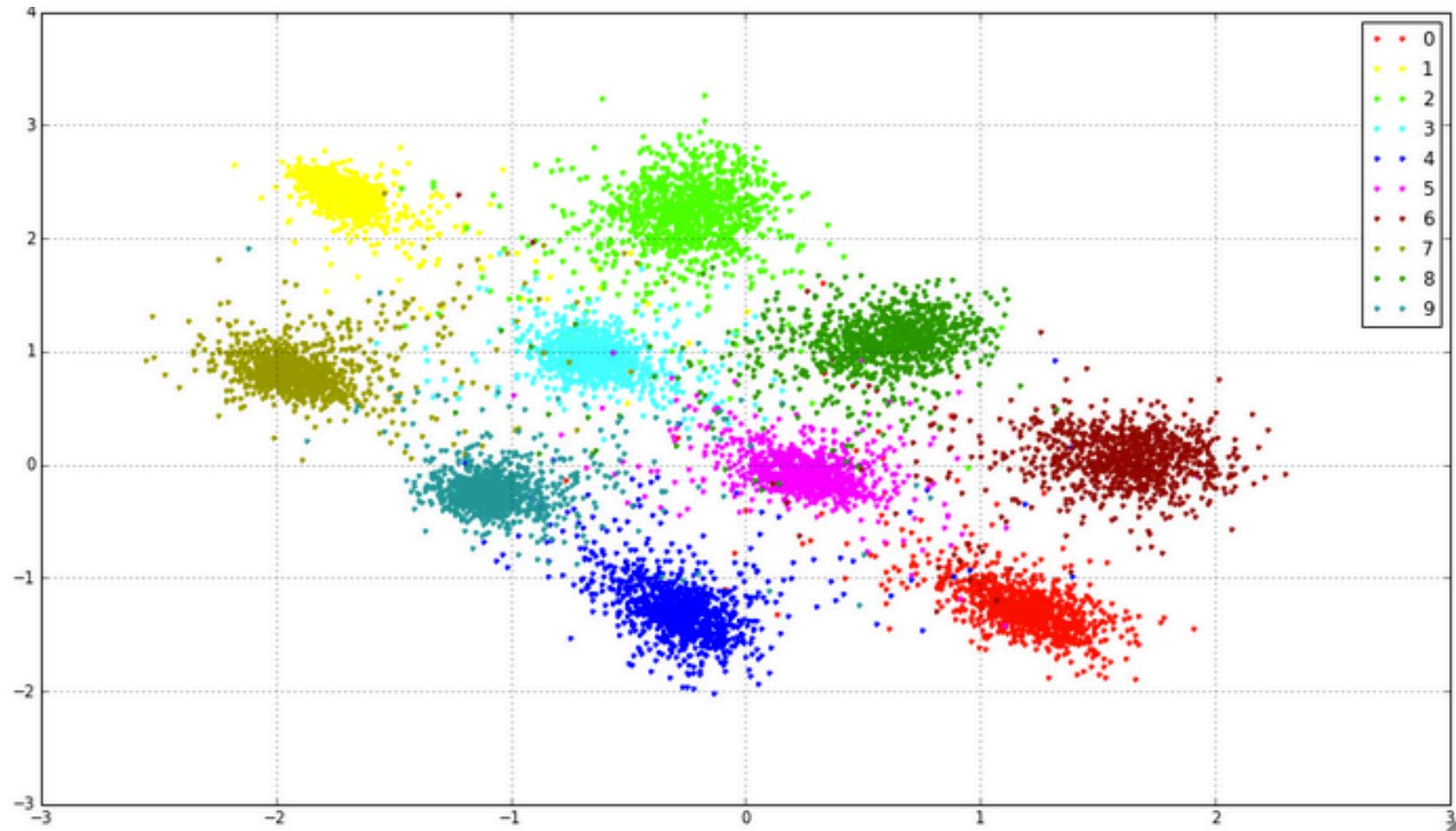
MNIST 500 iterations



MNIST 5000 iterations



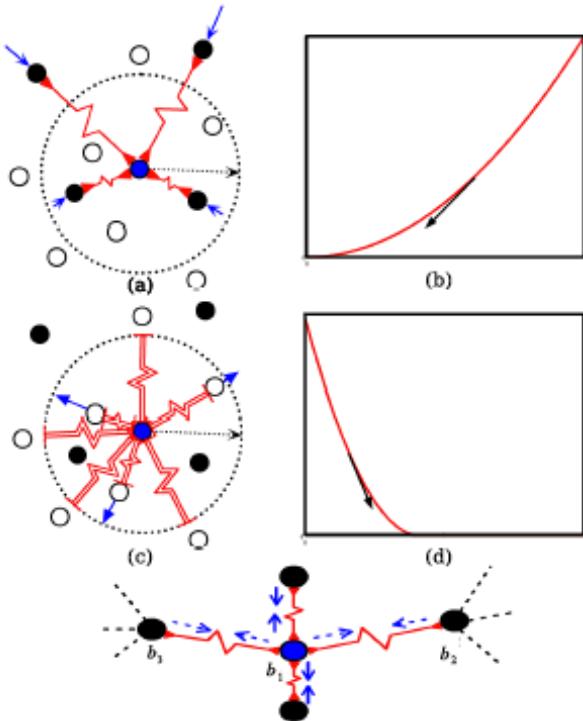
MNIST 50000 iterations



Advantages of Siamese Neural Networks

- Because you train on pairs of data points, they *amplify your training set*
- You don't need to know the number of categories in advance
- You don't even have to have examples for all categories!
- They are particularly suitable for domains where you have a large number of categories with a small number of examples

The loss function – The Spring Analogy



$$L(W, Y, \vec{X}_1, \vec{X}_2) =$$

$$(1 - Y) \frac{1}{2} (D_W)^2 + (Y) \frac{1}{2} \{ \max(0, m - D_W) \}^2$$

Where

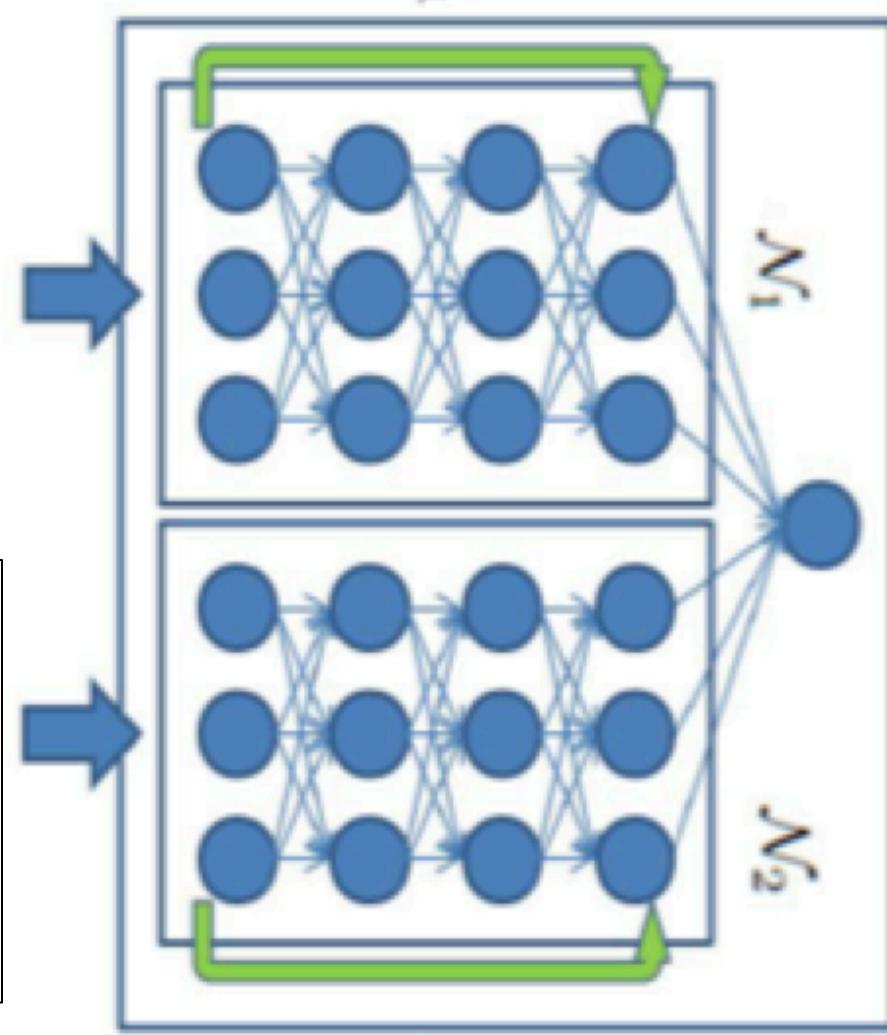
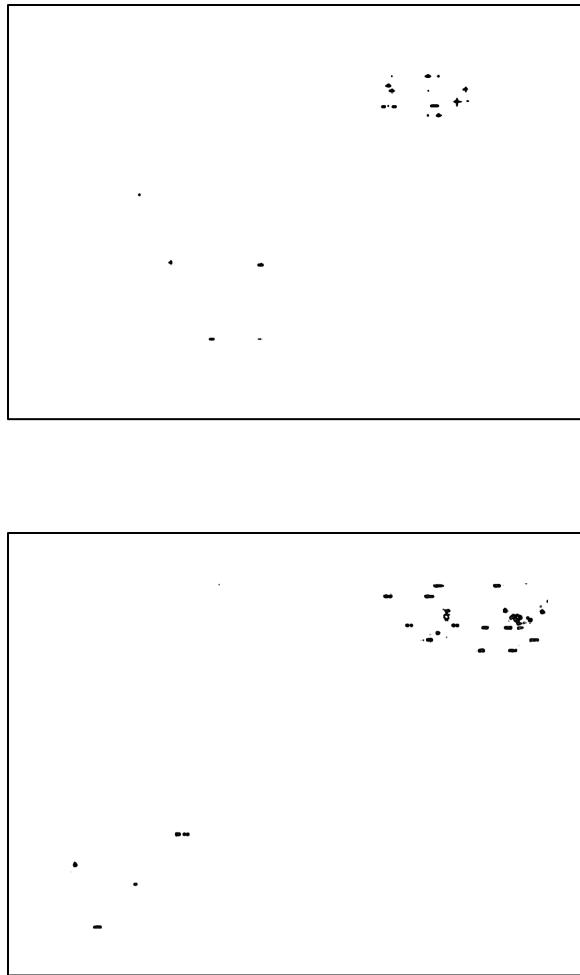
- D is the distance between the outputs of the two networks,
- m is a margin
- Y is 1 for a “different” pair and 0 for “same” pair

$$D_W(\vec{X}_1, \vec{X}_2) = \|G_W(\vec{X}_1) - G_W(\vec{X}_2)\|_2$$

$G_W(X)$ is the output vector of the network on stimulus X .

The basic idea

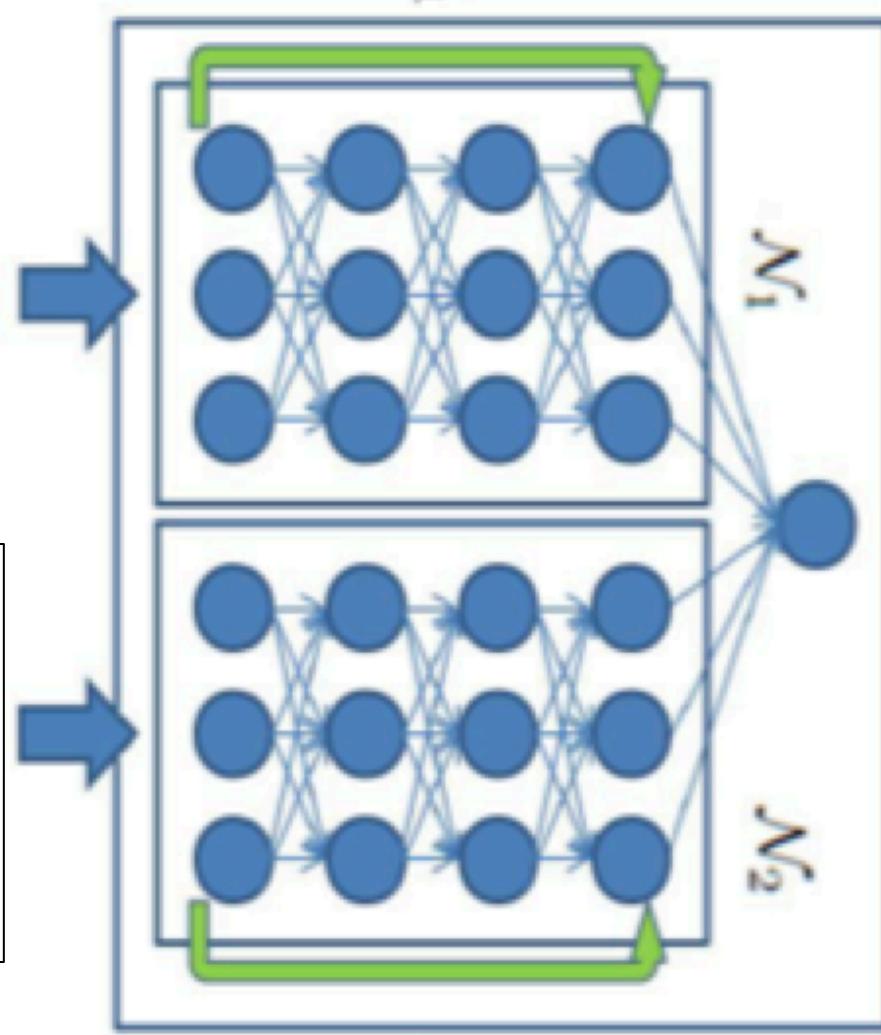
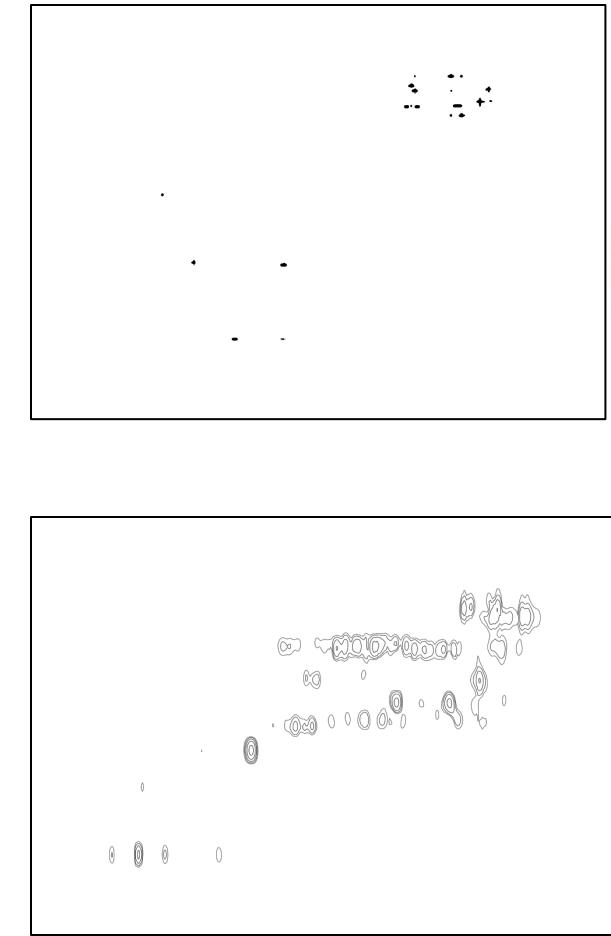
Aspewentin A Aspewentin C



“SAME”

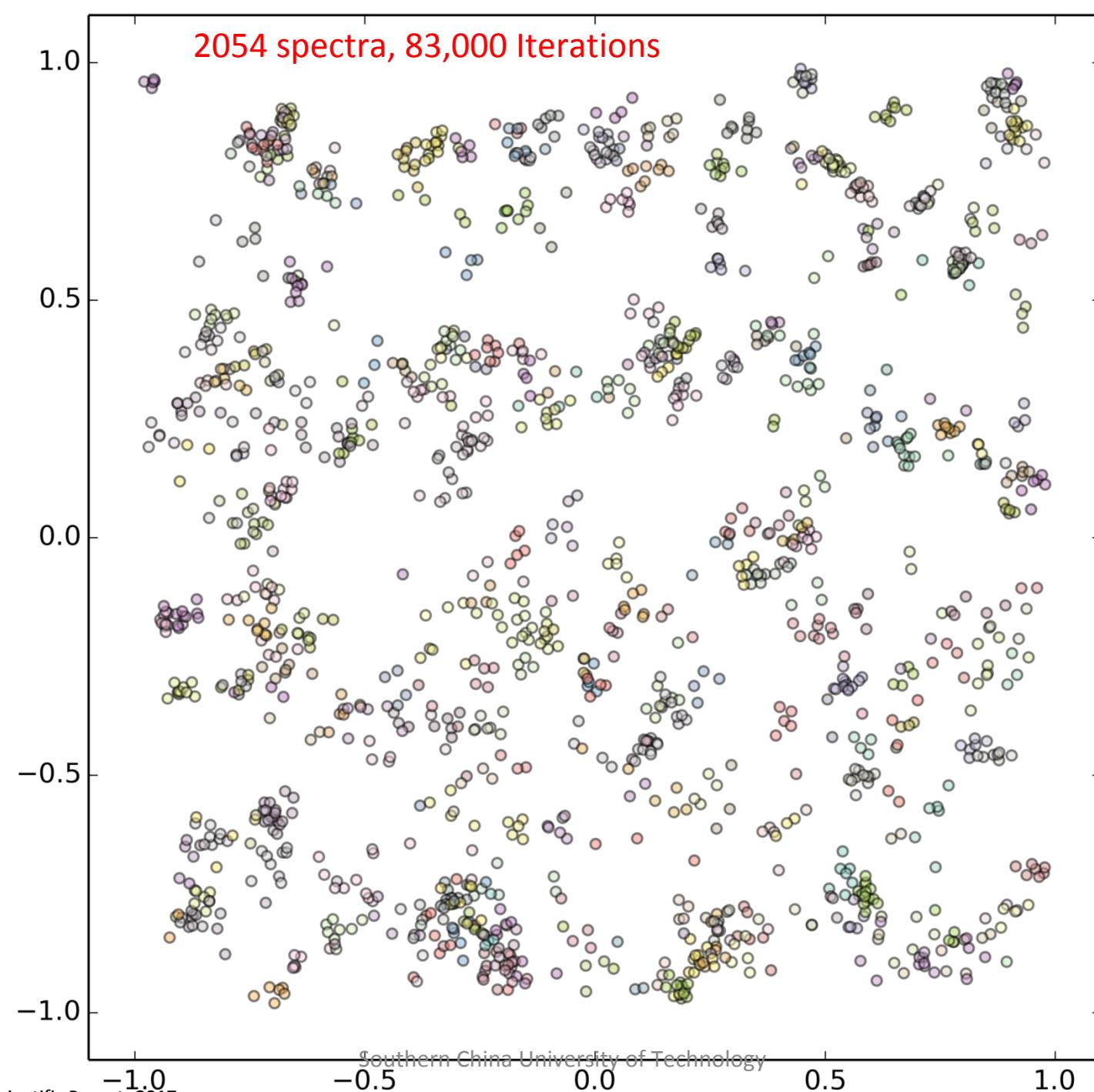
The basic idea

Aspewentin A Chandonanone A



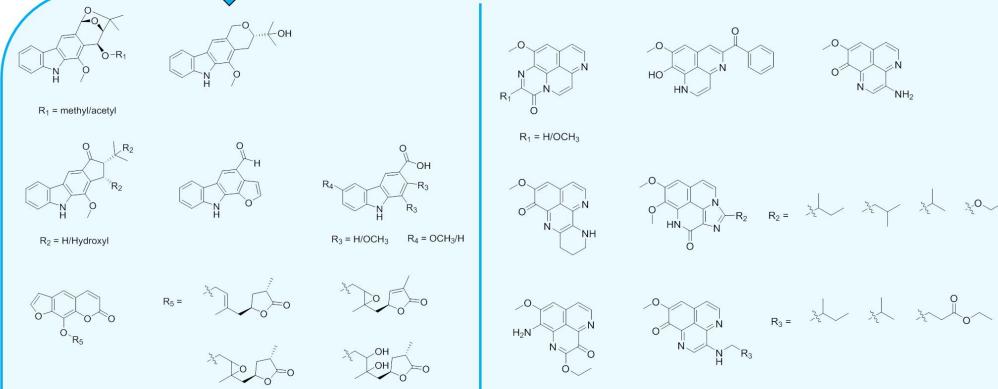
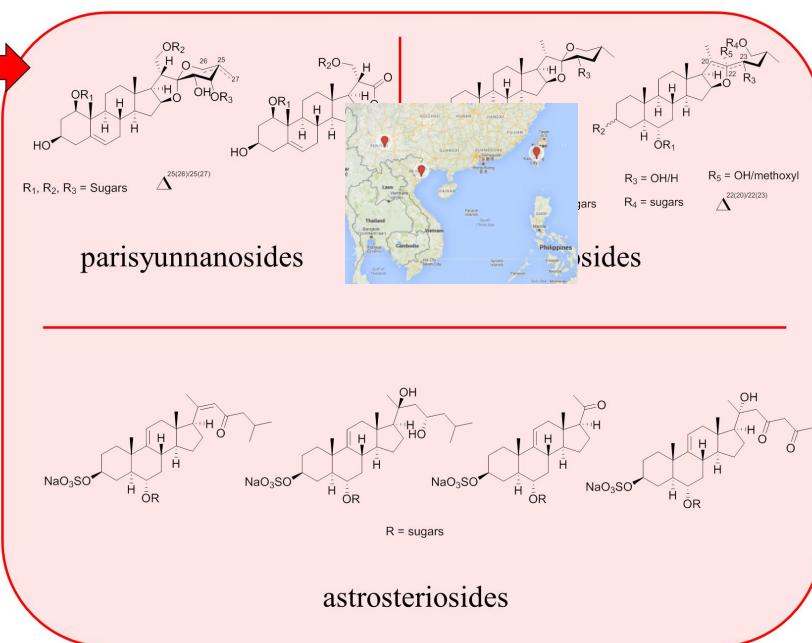
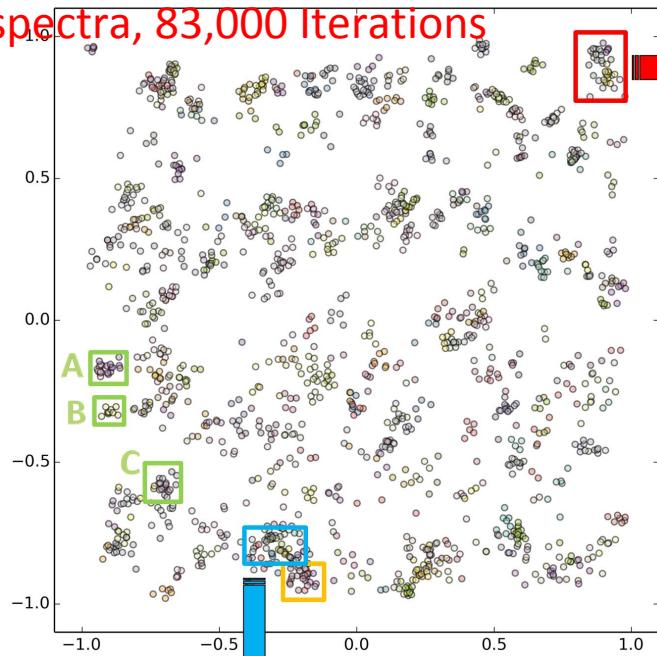
“DIFFERENT”

2054 spectra, 83,000 Iterations



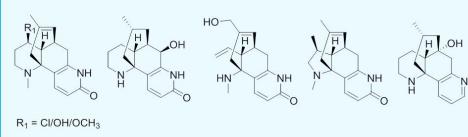
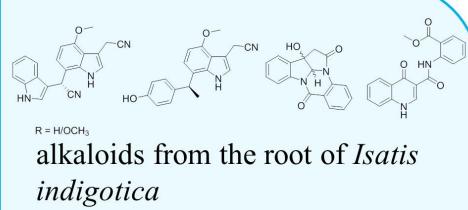
Qualitative Results

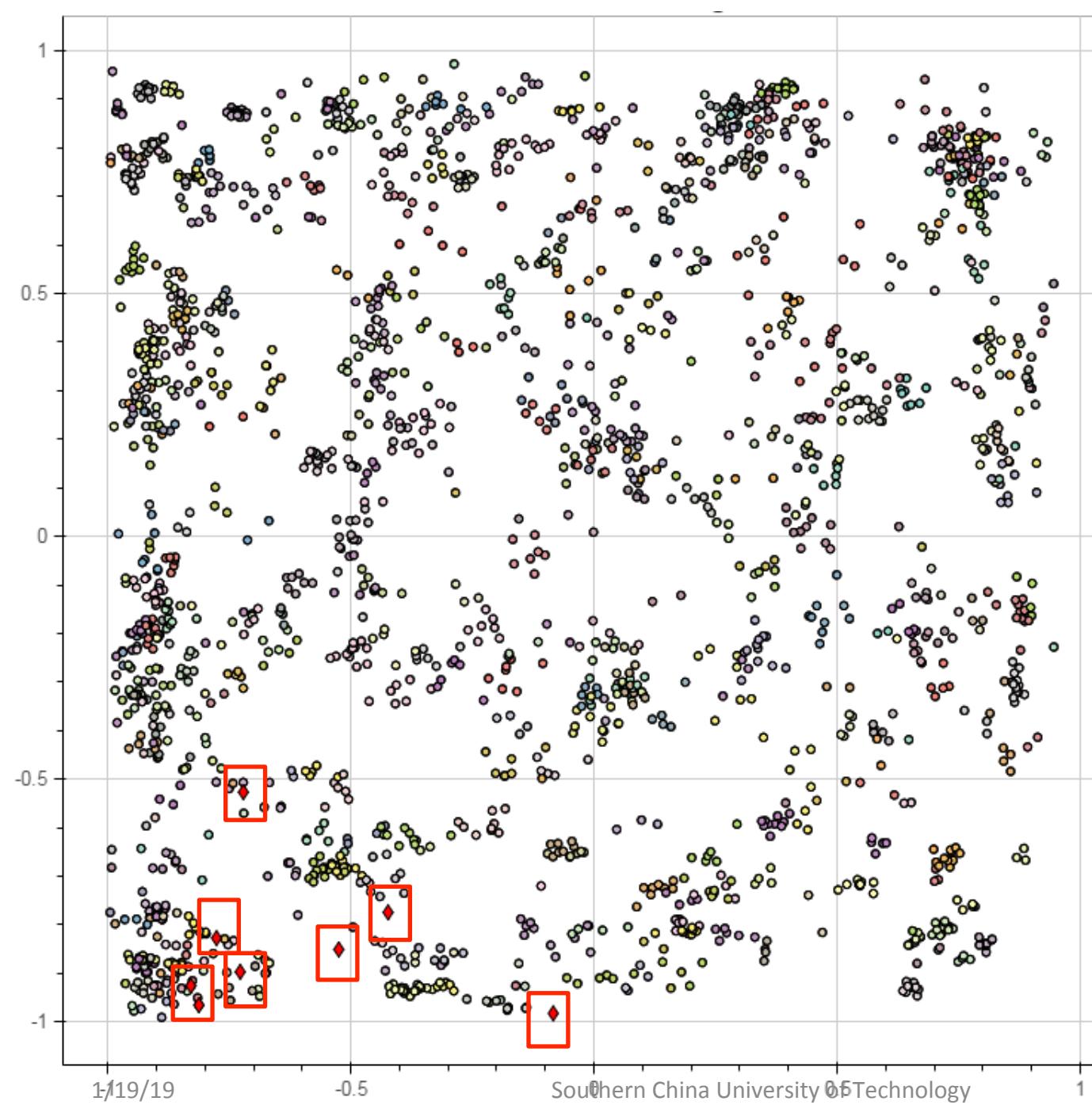
2054 spectra, 83,000 Iterations



roots of *Clausena lansium* derived compounds

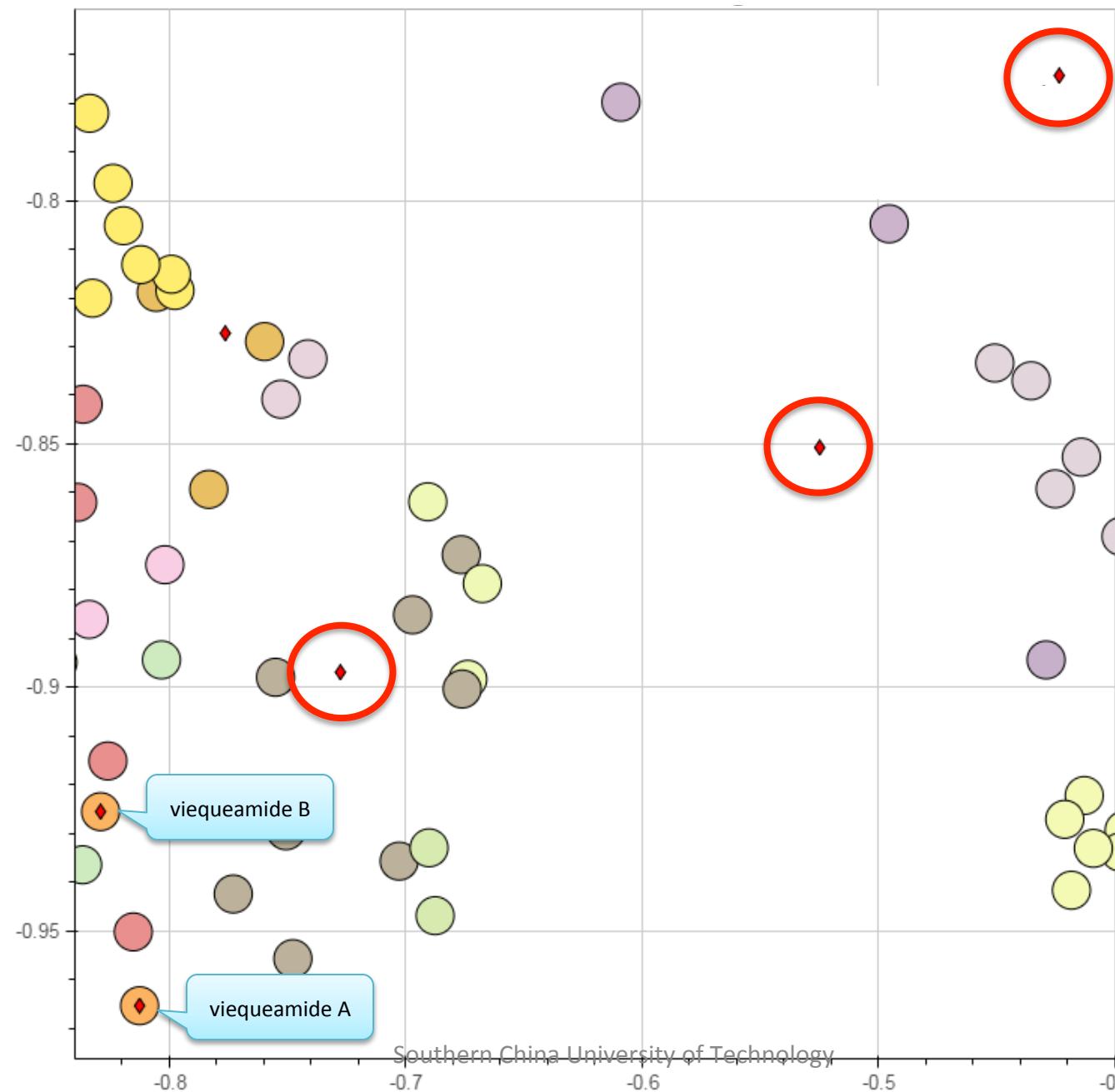
aaptamine derivatives from the Indonesian sponge *Aaptos suberitoides*, and the South China Sea sponge *Aaptos aaptos*



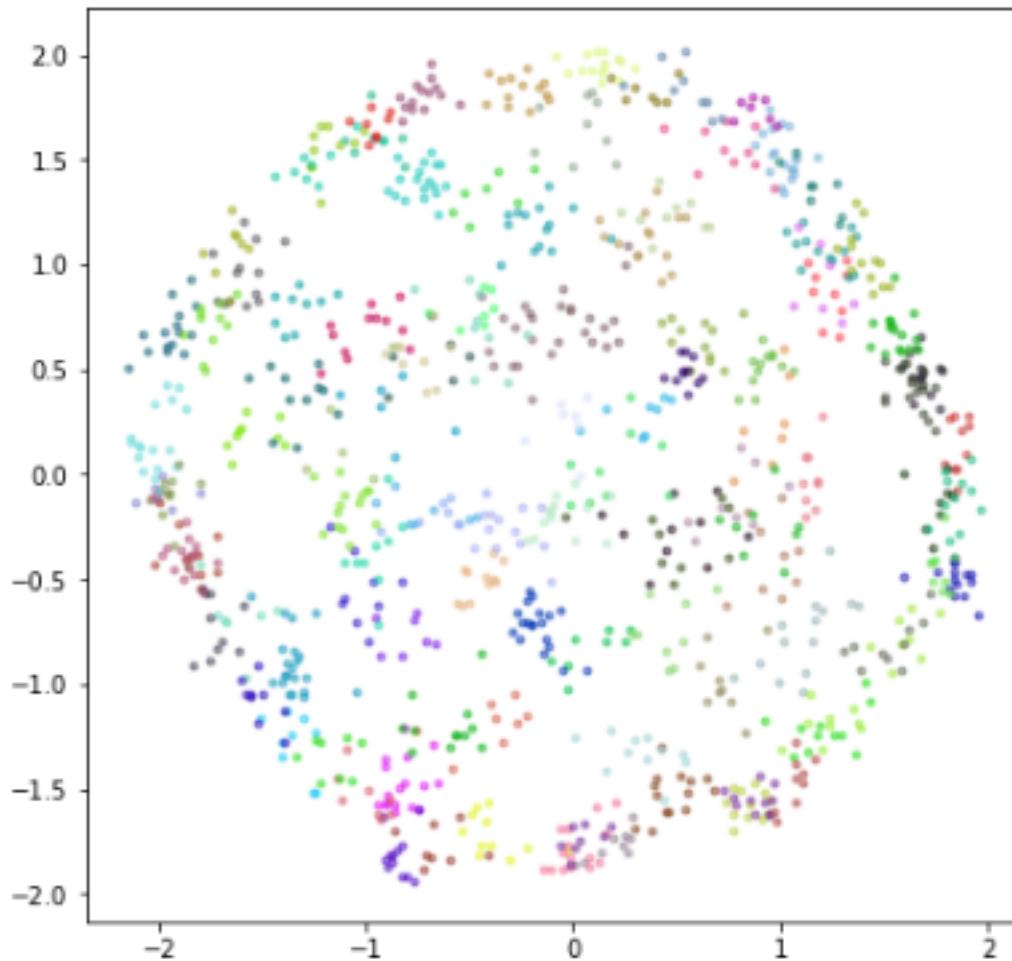


**Red Diamonds:
Embedding
of unknown
2D NMR
spectra into
the Cluster
Map**

Enlargement of the local area of the embedding map

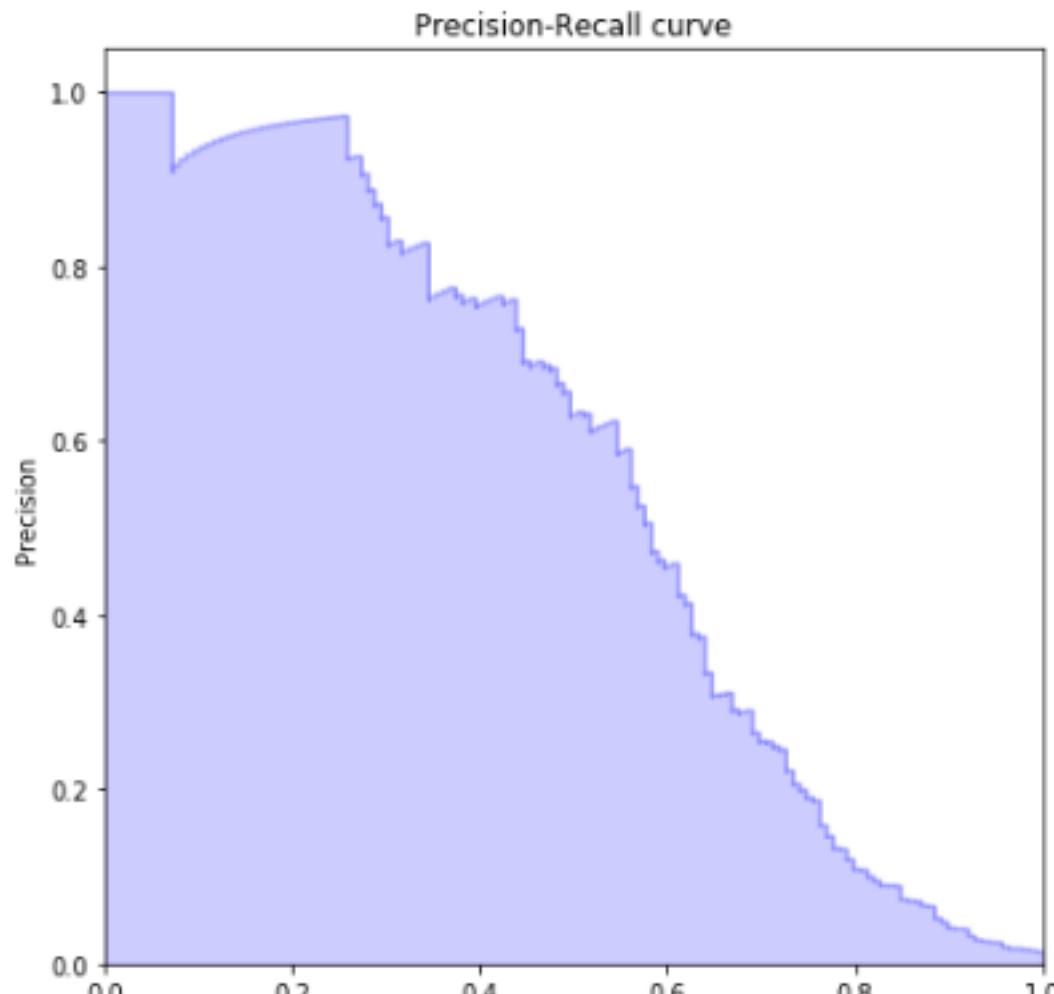


More recent results: just 300 epochs!



Quantitative Results: 10D outputs

Precision

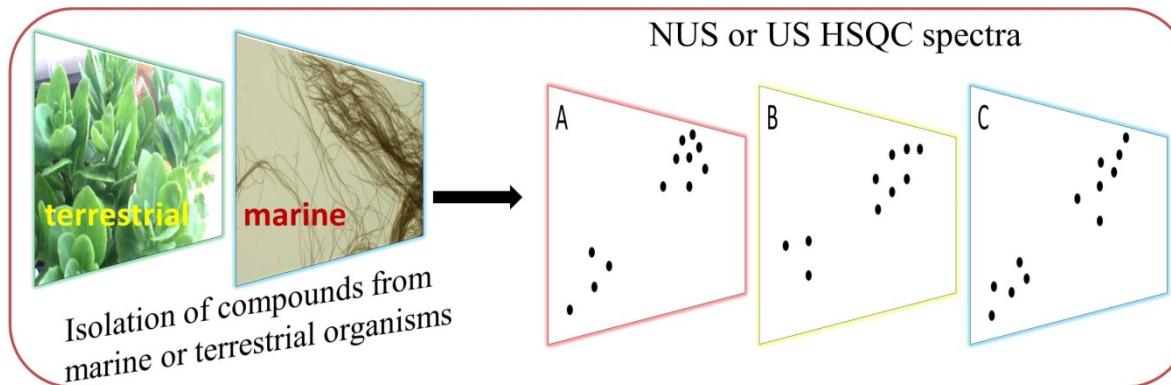


Recall

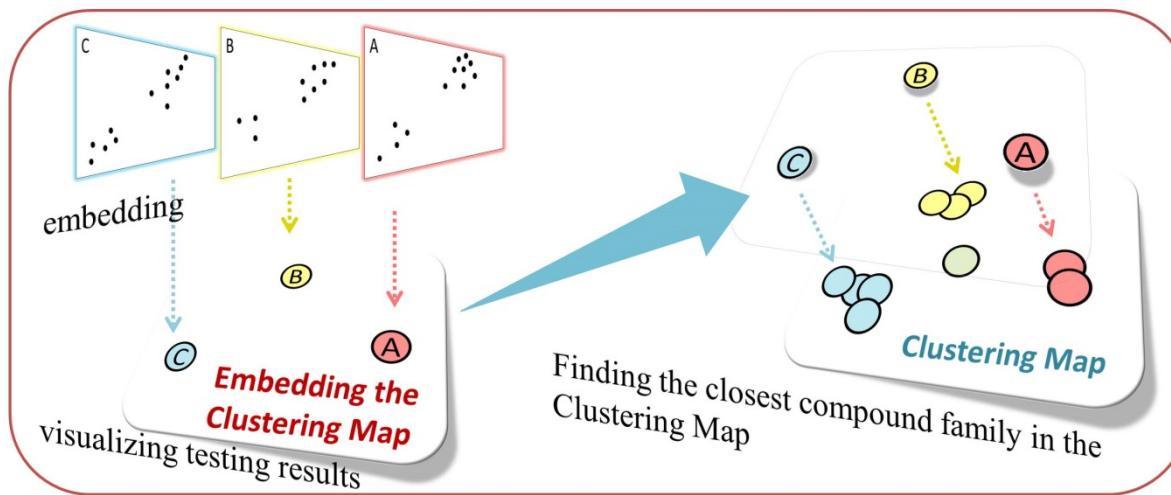
Demo

- What you will see is the result of a network with 3 outputs
- So, a 3D clustering

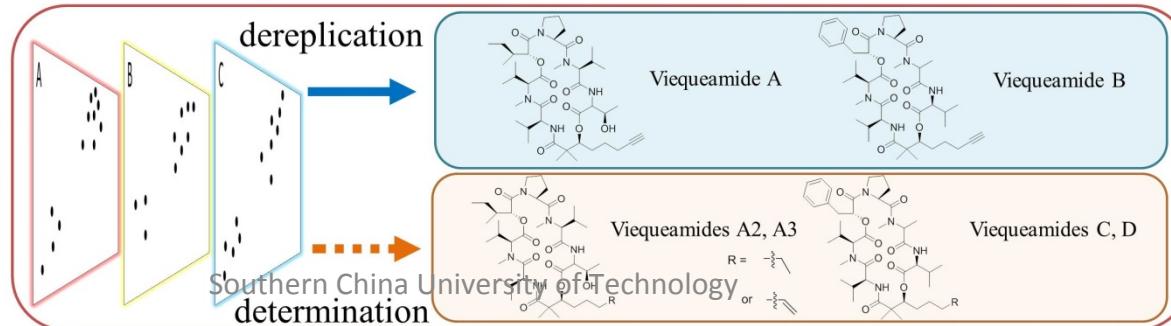
Data Collection



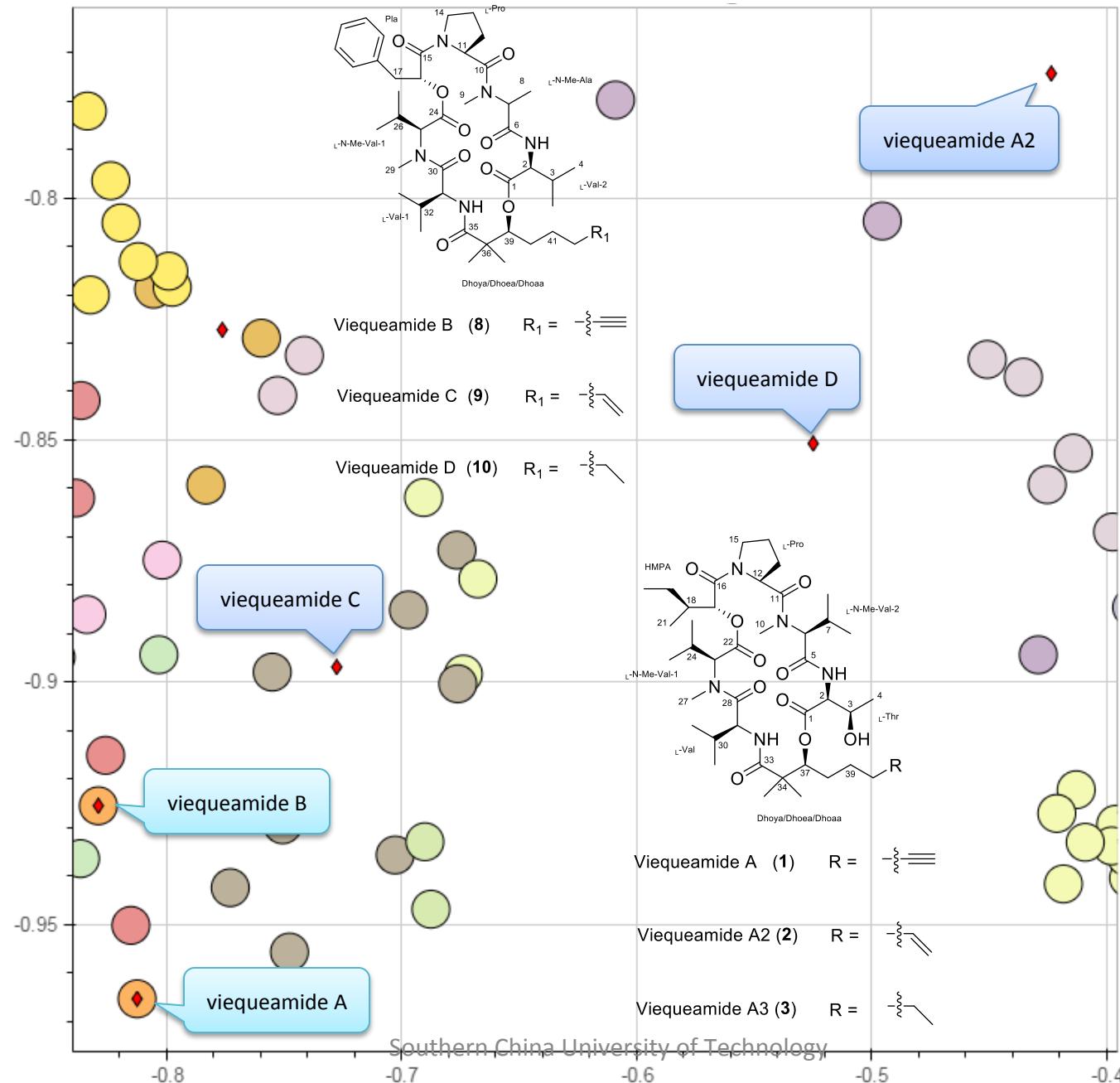
Data Analysis



Molecular Structures Output



Enlargement of the local area of the embedding map



Conclusions and Future Work

SMART Tutorial Upload About Donate Sign In

Coming Soon

SMART

Small Molecule Accurate Recognition Technology

EurekAlert! The Global Source for Science News AAAS

HOME NEWS

PUBLIC RELEASE: 7-11-17

SMART: structure

UNIVERSITY OF CA

Finding a Faster, More Accurate Way to Identify Molecular Structures of Natural Products

MEDIA CONTACT Doug Ramsey, (858) 522-5825.

An interdisciplinary team at the University of California San Diego has developed a method to identify the molecular structures of natural products that is significantly faster and more accurate than existing methods. The method

November 10, 2017

Roughly 70 percent of microorganisms in the soil and water have yet to be identified. The new technique can quickly analyze their molecular structures, helping researchers to better understand the environment. Materials science, Molecular structures, Synthetic biology, Environmental health, & Microbiology

Duke INTERDISCIPLINARY STUDIES INTERDISCIPLICINARITY BEYOND DUKE

Tag Archives: Facial recognition software

Week Ending Nov 3, 2017

SMART: Facial recognition for molecules

University of California San Diego researchers have developed a method to identify the molecular structures of natural products that is significantly faster and more accurate than existing methods. The method

GET SMART

MISS IT BY THAT MUCH!

- A lot of inquiries and requests from researchers worldwide
 - Pharmaceuticals
 - Perfume and cosmetics
 - Agricultural chemistry
 - Environmental sciences
 - NMR industries
- Soon: A user-friendly website with enhanced data security features

Future Work

- These results are using a “legacy” neural network – not optimized for the problem
 - Better network architecture
 - Using an autoencoding layer is supposed to improve results (Yann LeCun – personal communication)
- We are working on a new network objective function based on “purity” – a measure of cluster purity
 - Purity is non-differentiable
 - We have created a “soft”, differentiable version

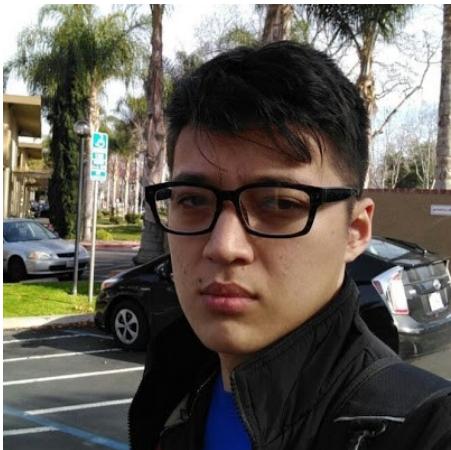
Acknowledgements

The NMR:

- Dr. Brendan Duggan
- Dr. Anthony Mrse
- Dr. Eugene Lin

Physics:

- Dr. Preston Landon
- Dr. Jie Min



Yerlan Idelbayev
UC Merced

1/19/19

GURU members:



Yash Nannapaneni



Nick Roberts



**Poornav
Sargur
Purushothama**

The Gerwick Lab



Southern China University of Technology

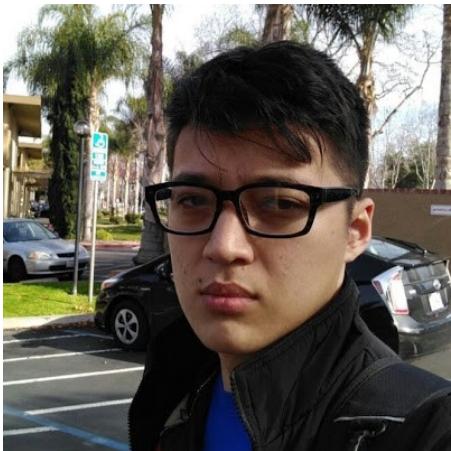
Acknowledgements

The NMR:

- Dr. Brendan Duggan
- Dr. Anthony Mrse
- Dr. Eugene Lin

Physics:

- Dr. Preston Landon
- Dr. Jie Min



Yerlan Idelbayev
UC Merced

1/19/19

GURU members:



Yash Nannapaneni



Nick Roberts



**Poornav
Sargur
Purushothama**

The Gerwick Lab



Bill Gerwick
Scripps Institution
of Oceanography
Southern China University of Technology



Chen Zhang
Scripps Institution
of Oceanography

Acknowledgements: Funders



Questions?