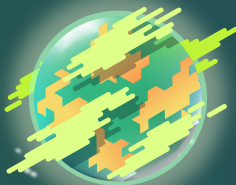# BEAM me up:
## A tale of Bounded Expansion Algorithms in Metagenomics

Blair D. Sullivan
NC State University
blair_sullivan@ncsu.edu
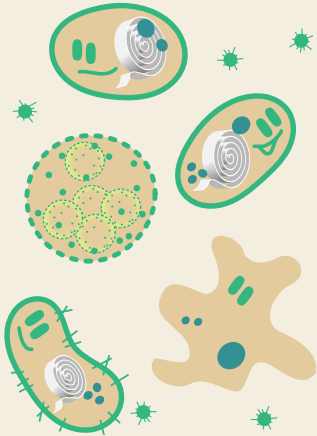
Workshop on Structural
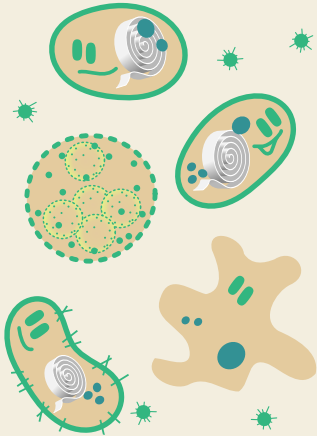Sparsity, Logic and Algorithms
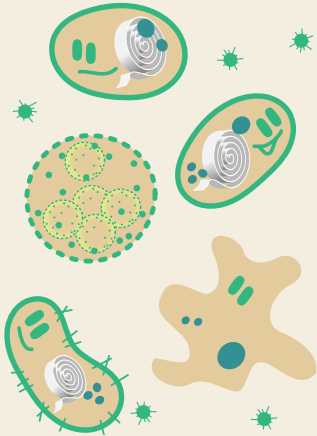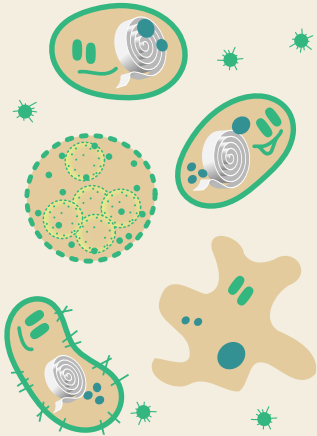
# Part I

# The scientific story

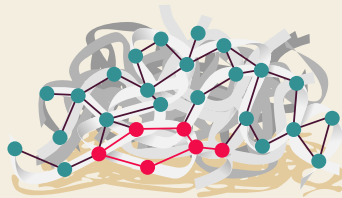# Metagenomics

# Metagenomics
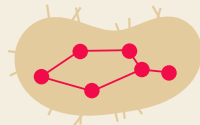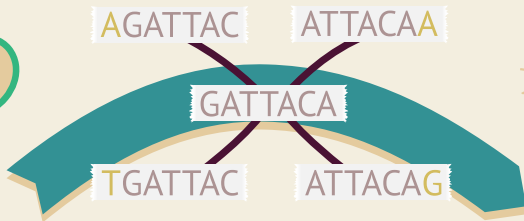
# Metagenomics

# De-Bruijn graphs

# De-Bruijn graphs

AGATTAC    ATTACAA

GATTACA

TGATTAC    ATTACAG

# De-Bruijn graphs

AGATTAC    ATTACAA

GATTACA

TGATTAC    ATTACAG

Bounded degree

# De-Bruijn graphs

AGATTAC    ATTACAA

GATTACA

TGATTAC    ATTACAG

Bounded degree

~100 million nodes

# CATLAS Overview

Genomic content

Domset

r-DTFAs

r-Domset

Dvořák's Algorithm*

de-Bruijn graphs

# Engineering Dvořák's algorithm



$$\mathrm{wcol}_{2r}$$

Approximation is
terrible in practice

$$\Delta^-(\vec{G}_{2r})$$

Approximation is
terrible in practice

Pseudocode   Implementable   Usable

Theory only   No tricks   Executable   Github
(or similar)

Dvořák Z. **Constant-factor approximation
of the domination number in sparse graphs**.
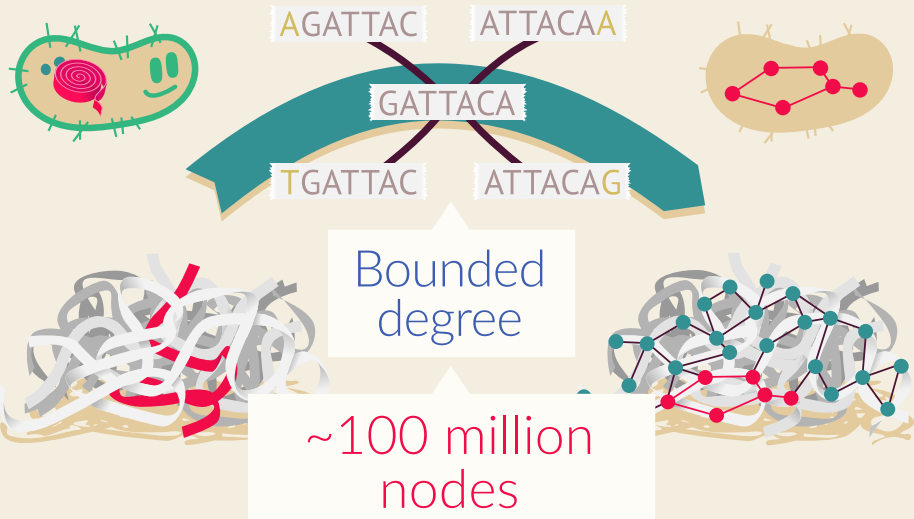European Journal of Combinatorics.
2013 Jul 31;34(5):833-40.

Reidl F. **Structural sparseness
and complex networks**.
(Doctoral dissertation, Dissertation,
Aachen, Techn. Hochsch., 2015).

# Engineering Dvořák's algorithm



$\mathrm{wcol}_{2r}$

Approximation is terrible in practise

$\Delta^-(\vec{G}_{2r})$

Approximation is terrible in practise

$\Delta^-(\vec{G}_r)$

Approximation is tunable (heuristic)

Pseudocode     Implementable     Usable

Theory only     No tricks     Executable     Github
(or similar)

Dvořák Z. **Constant-factor approximation of the domination number in sparse graphs.** European Journal of Combinatorics. 2013 Jul 31;34(5):833-40.
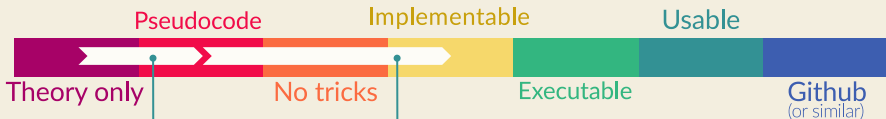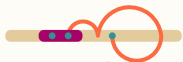
**Barnraising for Data-Intensive Discovery** at MDI Biological Laboratory in Salisbury Cove, Maine.

Reidl F. **Structural sparseness and complex networks.** (Doctoral dissertation, Dissertation, Aachen, Techn. Hochsch., 2015).

Ongoing collaboration with **Theory in Practice Group** (NCSU) and **Lab for Data Intensive Biology** (UC Davis) https://github.com/spacegraphcats/spacegraphcats

# Engineering Dvořák's algorithm

For input a graph G and integers r,t our algorithm computes an $2(t+2)\Delta^-(\vec{G}_{2r})\Delta^-(\vec{G}_r)$ -approximate r-dominating set. Importantly, it computes only the rth dtf-augmentation $\vec{G}_r$.

Fudge-factor t: small t yield better approximation guarantee in the worst case, but larger dominating sets in practice!

Pseudocode    Implementable    Usable

Theory only    No tricks    Executable    Github
(or similar)

Dvořák Z. **Constant-factor approximation of the domination number in sparse graphs**. European Journal of Combinatorics. 2013 Jul 31;34(5):833-40.

Reidl F. **Structural sparseness and complex networks**. (Doctoral dissertation, Dissertation, Aachen, Techn. Hochsch., 2015).
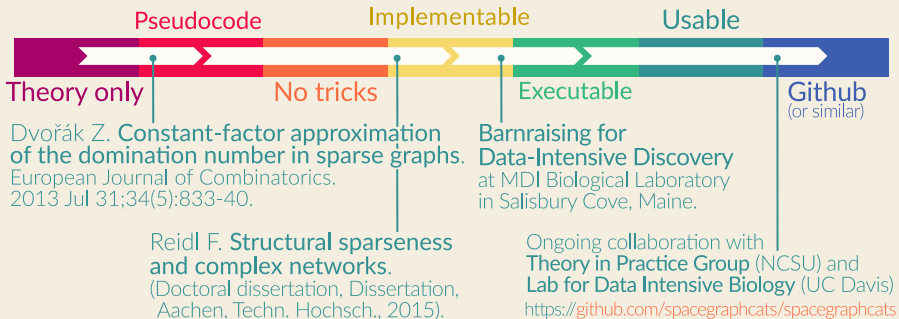
**Barnraising for Data-Intensive Discovery** at MDI Biological Laboratory in Salisbury Cove, Maine.

Ongoing collaboration with **Theory in Practice Group** (NCSU) and **Lab for Data Intensive Biology** (UC Davis) https://github.com/spacegraphcats/spacegraphcats

CATLAS-1 Results

CATLAS-1 Results

CATLAS-1 Results

Kingdom
- Archaea
- Bacteria

**Total # Occurrences of any Ortholog**

KEGG Ortholog BRITE Hierarchy

Biosynthesis of amino acids
Carbon metabolism
Ribosome
Methane metabolism
Purine metabolism
Pyrimidine metabolism
Aminoacyl–tRNA biosynthesis
2–Oxocarboxylic acid metabolism
ABC transporters
Carbon fixation pathways in prokaryotes
Pyruvate metabolism
Glycolysis / Gluconeogenesis
Cysteine and methionine metabolism
Valine, leucine and isoleucine biosynthesis
Glycine, serine and threonine metabolism
Amino sugar and nucleotide sugar metabolism
Alanine, aspartate and glutamate metabolism
Quorum sensing
Pentose phosphate pathway
RNA degradation
Citrate cycle (TCA cycle)
C5–Branched dibasic acid metabolism
Oxidative phosphorylation
Lysine biosynthesis
Butanoate metabolism
Fructose and mannose metabolism
Phenylalanine, tyrosine and tryptophan biosynthesis
Histidine metabolism
Arginine biosynthesis
Propanoate metabolism
Folate biosynthesis
Porphyrin and chlorophyll metabolism
Pantothenate and CoA biosynthesis
Glyoxylate and dicarboxylate metabolism
DNA replication
Peptidoglycan biosynthesis
Mismatch repair
Homologous recombination
Carbon fixation in photosynthetic organisms
RNA polymerase
Protein export
Nucleotide excision repair
Cell cycle – Caulobacter
One carbon pool by folate
Starch and sucrose metabolism
Base excision repair
Nicotinate and nicotinamide metabolism
Bacterial secretion system
Drug metabolism – other enzymes
Arginine and proline metabolism
Thiamine metabolism
Terpenoid backbone biosynthesis
Sulfur relay system
Streptomycin biosynthesis
Selenocompound metabolism
Fatty acid biosynthesis
Bacterial chemotaxis
Vancomycin resistance
Glycerophospholipid metabolism
Valine, leucine and isoleucine degradation
Riboflavin metabolism
Nitrotoluene degradation
beta–Alanine metabolism
Biotin metabolism
Ubiquinone and other terpenoid–quinone biosynthesis
Lipopolysaccharide biosynthesis
Taurine and hypotaurine metabolism

count
0      100      200      300      400

A long time ago, in a galaxy far, far away...

"Sparse Graph Cuts"

Spacegraphcats

Genomic content

Domset

r-DTFAs

r-Domset

Dvořák's Algorithm

GORDON AND BETTY
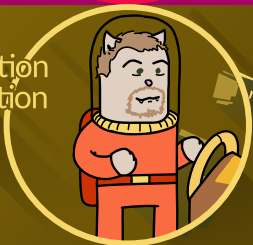MOORE
FOUNDATION

# slack

HackMD

GitHub

jupyter

18 Months later.

Translation
Method design
Course planning

Data-wrangling
Target identification
Biological validation
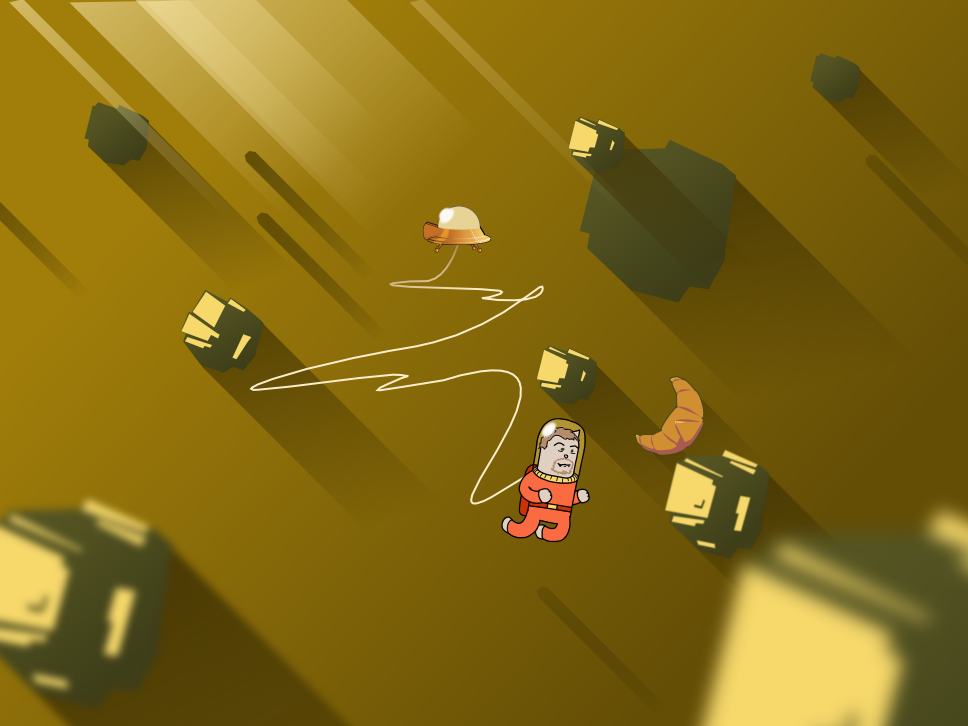
Scalability
Feature development
Implementation
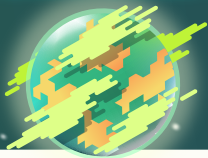
The efficient computation of dominating sets will open up a whole new range of possibilities in biology.
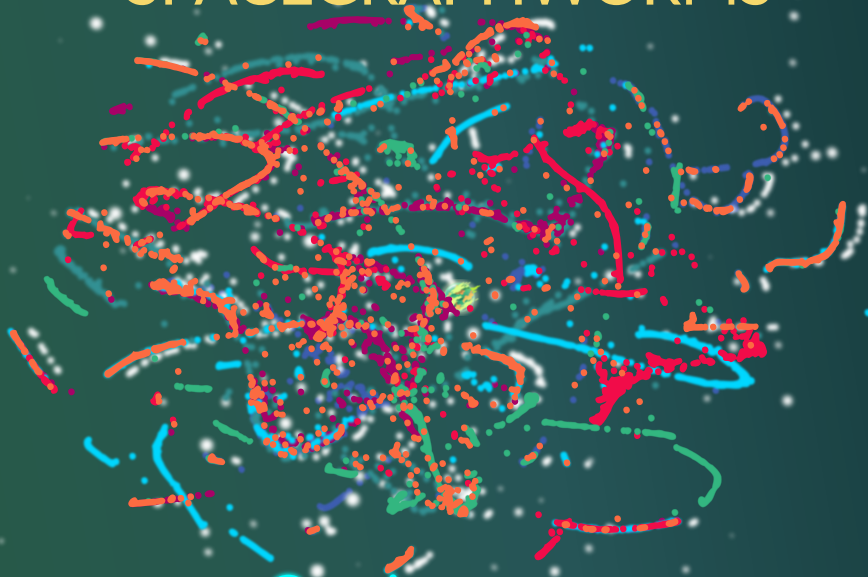
C. Titus Brown, professor at UC Davis

spacegraphcats will transform the way biologists interact with genome assemblies.

It allows us to access previously discarded sequencing information thereby allowing more robust functional characterization.

Taylor Reiter, his much more eloquent student

# SPACEGRAPHWORMS

Coming soon...