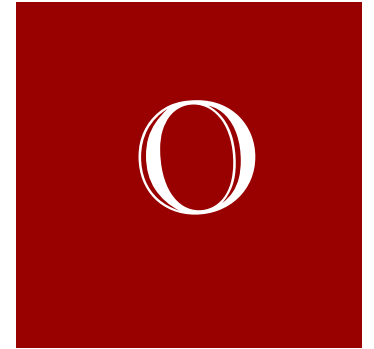




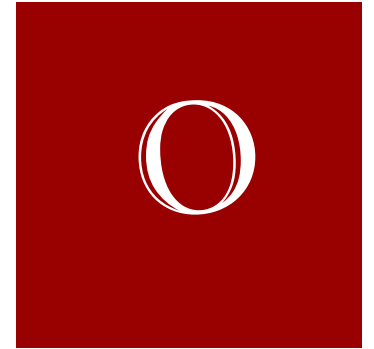
Algorithms and Methods for Cross-Species Forensics and Classification

Darryl Reeves, Ph.D.
Clinical and Research Genomics
04/11/2018



Outline

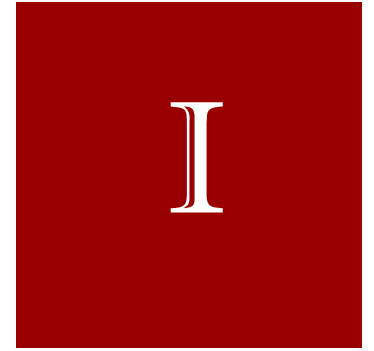
- I. Introduction
- II. An alternative view of biological sequences
- III. A model for sequence identification
- IV. Conclusion/Future plans



Outline

- I. Introduction
- II. An alternative view of biological sequences
- III. A model for sequence identification
- IV. Conclusion/Future Research Agenda

Why study communities of organisms (metagenomics)?




- Environment contains some interesting biological activity
- Interaction between community and external environment/host
- Understand how environmental changes impact organisms

Motivating examples

Metagenome, metatranscriptome and single-cell sequencing reveal microbial response to Deepwater Horizon oil spill

Olivia U Mason^{1,2}, Terry C Hazen^{1,3}, Sharon Borglin¹, Patrick SG Chain^{4,5}, Eric A Dubinsky¹, Julian L Fortney¹, James Han^{4,5}, Hoi-Ying N Holman¹, Jenni Hultman¹, Regina Lamendella¹, Rachel Mackelprang⁵, Stephanie Malfatti^{5,6}, Lauren M Tom¹, Susannah G Tringe⁵, Tanja Woyke⁵, Jizhong Zhou^{7,8}, Edward M Rubin⁵ and Janet K Jansson^{1,5}

The Dynamics of the Human Infant Gut Microbiome in Development and in Progression toward Type 1 Diabetes

Aleksandar D. Kostic, Dirk Gevers, Heli Siljander, Tommi Vatanen, Tuulia Hyötyläinen, Anu-Maaria Hämäläinen, Aleksandr Peet, Vallo Tillmann, Päivi Pöhö, Ismo Mattila, Harri Lähdesmäki, Eric A. Franzosa, Outi Vaarala, Marcus de Goffau, Hermie Harmsen, Jorma Ilonen, Suvi M. Virtanen, Clary B. Clish, Matej Orešič, Curtis Huttenhower, Mikael Knip²³ on behalf of the DIABIMMUNE Study Group²², Ramnik J. Xavier²³  

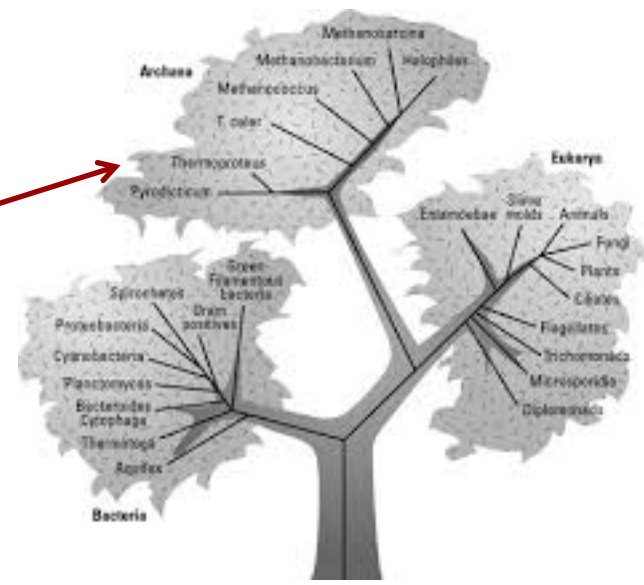
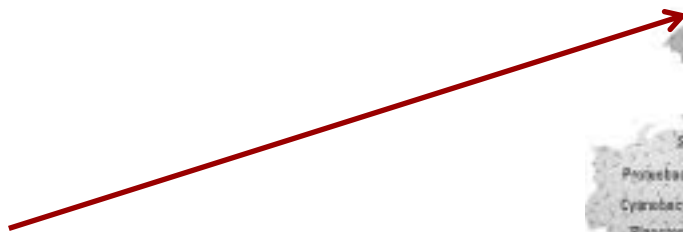
Metagenomics research

- Two very different research domains (environmental and medical)
- Same basic questions
- Given a heterogeneous community of organisms
 - Who is there?
 - What are they doing (function)?

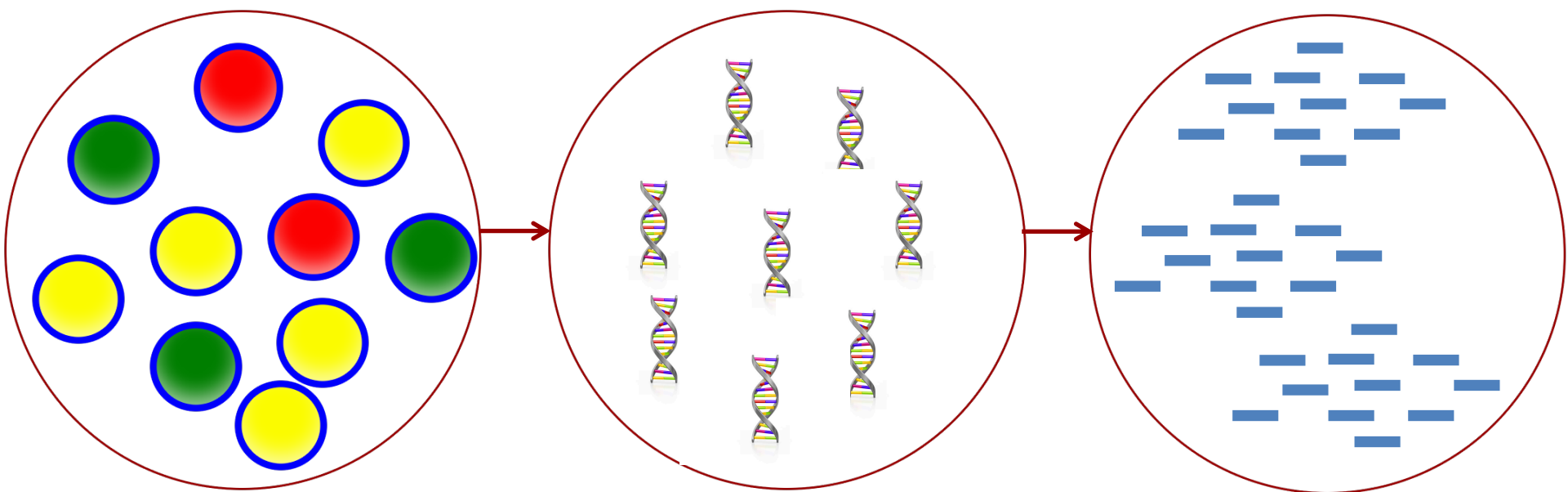
The Goal



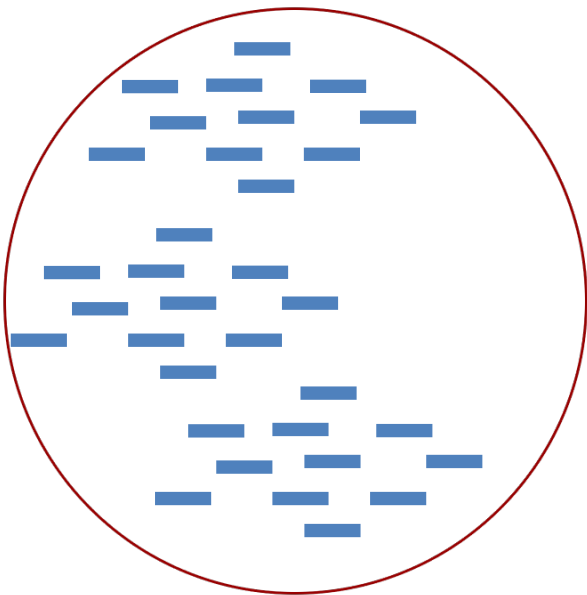
DNA Sequence



Metagenomics – data



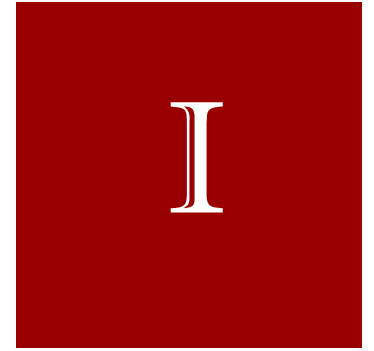
Metagenomics – data



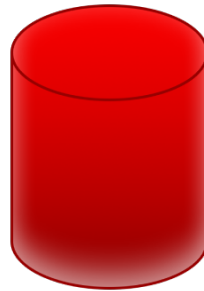
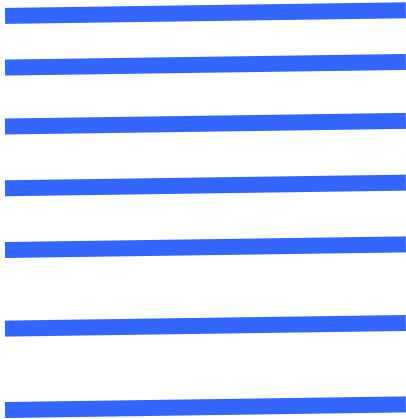
Metagenomics – data



Metagenomic identification/classification



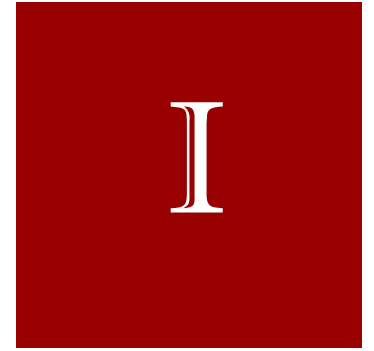
Species X



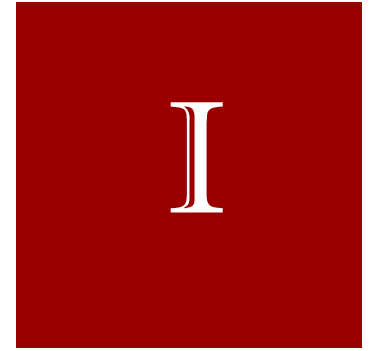
- Sequence alignment
- Sub-sequence markers
- Sub-sequence composition
- Hybrid approaches

Common Limitation

Metagenomic identification/classification



Metagenomic identification/classification



Species X

Species Y

Species Z

Species 1

Species 2

Species 3

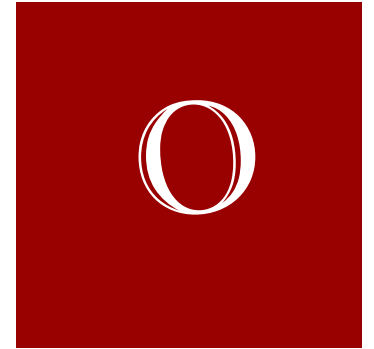


Hypothesis

Current methods fail to utilize all available information

Results in a local view of community composition

Statistical learning can be used to develop a global view of communities



Outline

- I. Introduction
- II. An alternative view of biological sequences
- III. A model for sequence identification
- IV. Conclusion/Future plans



An alternative view

- Sequence identification fundamentally a comparison problem

- Comparing text vs. numbers



K-mers

ATGCTAGGAACCCTAGCTTACAGAGCAGTTGCAGG



K-mers

ATGCTAGGAACCCTAGCTTACAGAGCAGTTGCAGG



K-mers

ATGCTAGGAACCCTAGCTTACAGAGCAGTTGCAGG



K-mers

ATGCTAGGAACCCTAGCTTACAGAGCAGTTGCAGG



K-mers

ATGCTAGGAACCCTAGCTTACAGAGCAGTTGCAGG



K-mers

ATGCTAGGAACCCTAGCTTACAGAGCAGTTGCAGG



K-mers

ATGCTAGGAACCCTAGCTTACAGAGCAGTTGCAGG



K-mers

ATGCTAGGAACCCTAGCTTACAGAGCAGTTGCAGG

3-mer	Count
ATG	1
TGC	2
CTA	2
TAG	2
AGG	2
GGA	1
GAA	1

K-mer Frequency Distribution



K-mers

ATGCTAGGAACCCTAGCTTACAGAGCAGTTGCAGG

3-mers, 64 combinations

5-mers, 1024 combinations

23-mers, 70,368,744,177,664 combinations

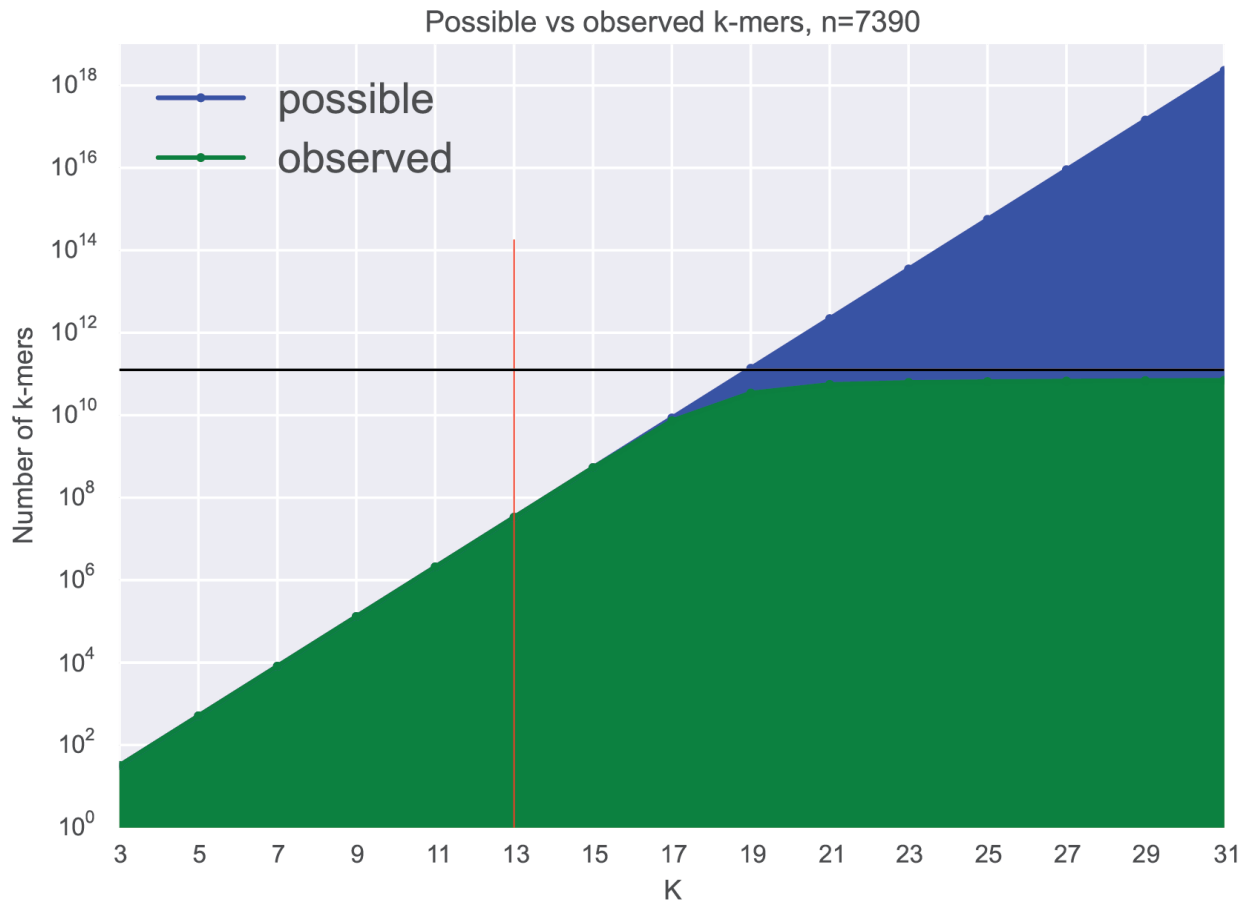
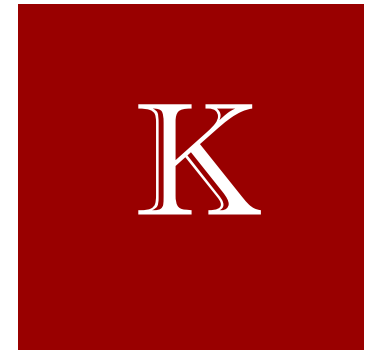


K-mers

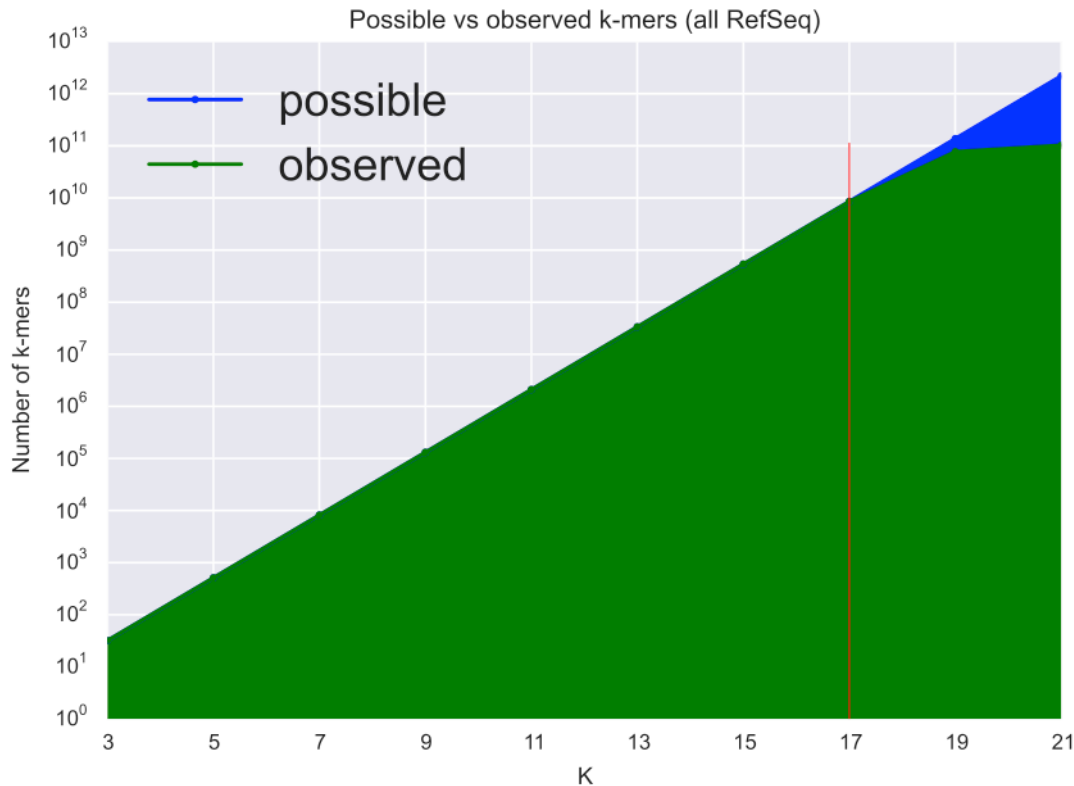
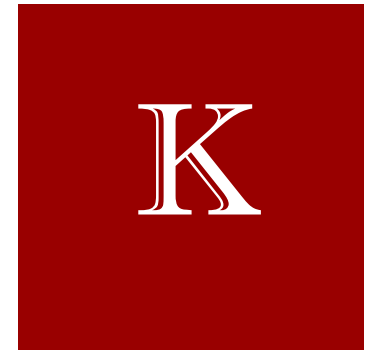
ATGCTAGGAACCCTAGCTTACAGAGCAGTTGCAGG

4^K

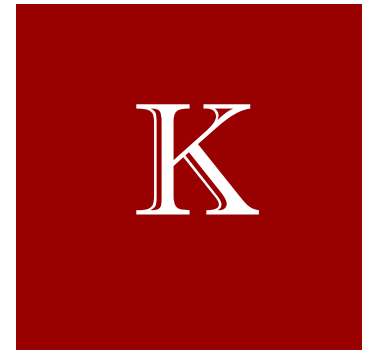
K-mer diversity analysis – excluding k-mers



K-mer diversity analysis – excluding k-mers



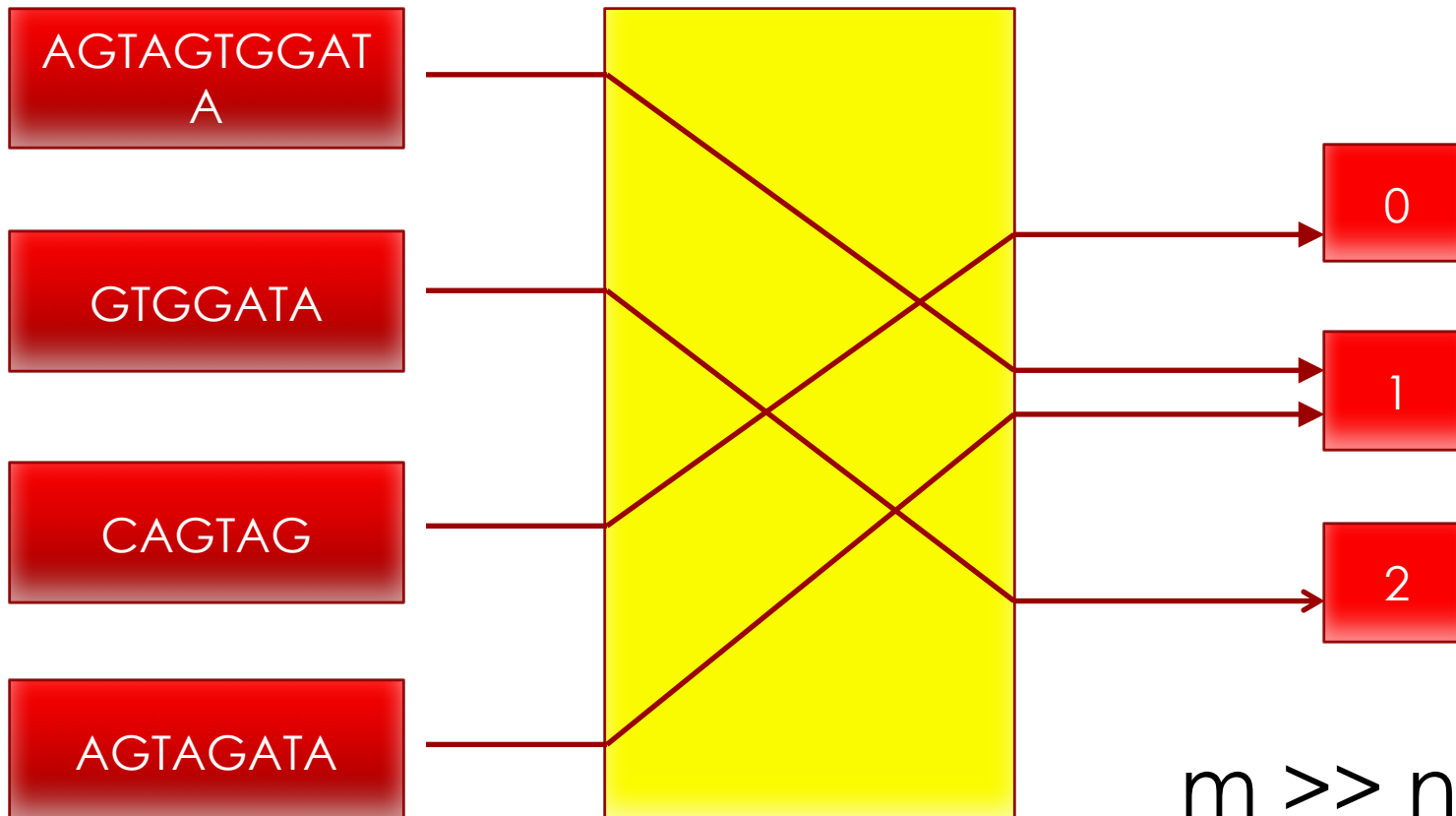
Hashing



Keys (m)

Hash Function

Hashes (n)

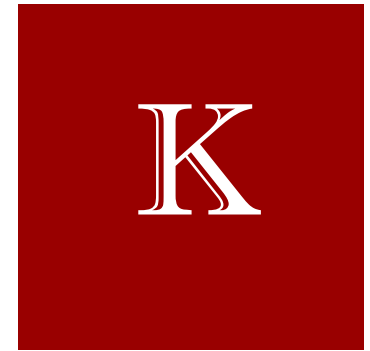


Data from KMerge



Sequence S_i	
Hash Val	Count
0	4
1	0
2	9
\vdots	\vdots
$ H -3$	13
$ H -2$	1
$ H -1$	2

KMerge - reference



FASTA

```
>genome1
AGCTCTATCATGCCCTTAGTT
TTAAACTAGGTCTAAGCTAG
AAGG...
```

+

Taxonomy

Kingdom	x
Domain	xy
Phylum	xz
Class	xa
Order	xb
Family	xt
Genus	xx
Species	y



Count K-mers

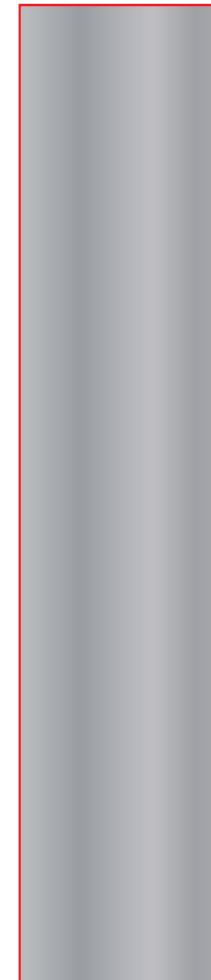
GGGCATCTAT12	
TAATC	65
AGT	2
ACCTA	1
ACGATCATGC78	
ACGTT	921
...	
CAT	2702
TGCCTGA	53
CTGAT	12
GCA	97
GCATA	231
TTATATC	780

Hash K-mers

123214	11
6350840	310
534543535	391
1000001231	1320
1232144523	32434
1432234243	8532
2301430583	43240
2434238402	545
3424320890	95435
3432800243	14234
4124023804	89243



DB



FASTA

```
>genome2
AGCTCTATCATGCCCTTAGTT
TTAAACTAGGTCTAAGCTAG
AAGG...
```

+

Taxonomy

Kingdom	z
Domain	zy
Phylum	zz
Class	za
Order	zb
Family	zt
Genus	zx
Species	b



GCTCAAATAT12	
TAATCGG	22
GGTAAG	5
ACCTA	45
ACGATCATGC70	
ACG	1901
...	
CATGG	702
TGTGA	59
CTGATAC	14
GCAGTCA	209
ATA	210
TTAGGTA	283

12321114	23
5033840	652
55351545	31
100231	232
123223807	3241
1432243	8582
2303058331	43012
2438402312	567
6420890	9551
3400243537	1342
402380443	8431



FASTA

```
>genome3
AGCTCTATCATGCCCTTAGTT
TTAAACTAGGTCTAAGCTAG
AAGG...
```

+

Taxonomy

Kingdom	w
Domain	wy
Phylum	wz
Class	wa
Order	wb
Family	wt
Genus	wx
Species	d

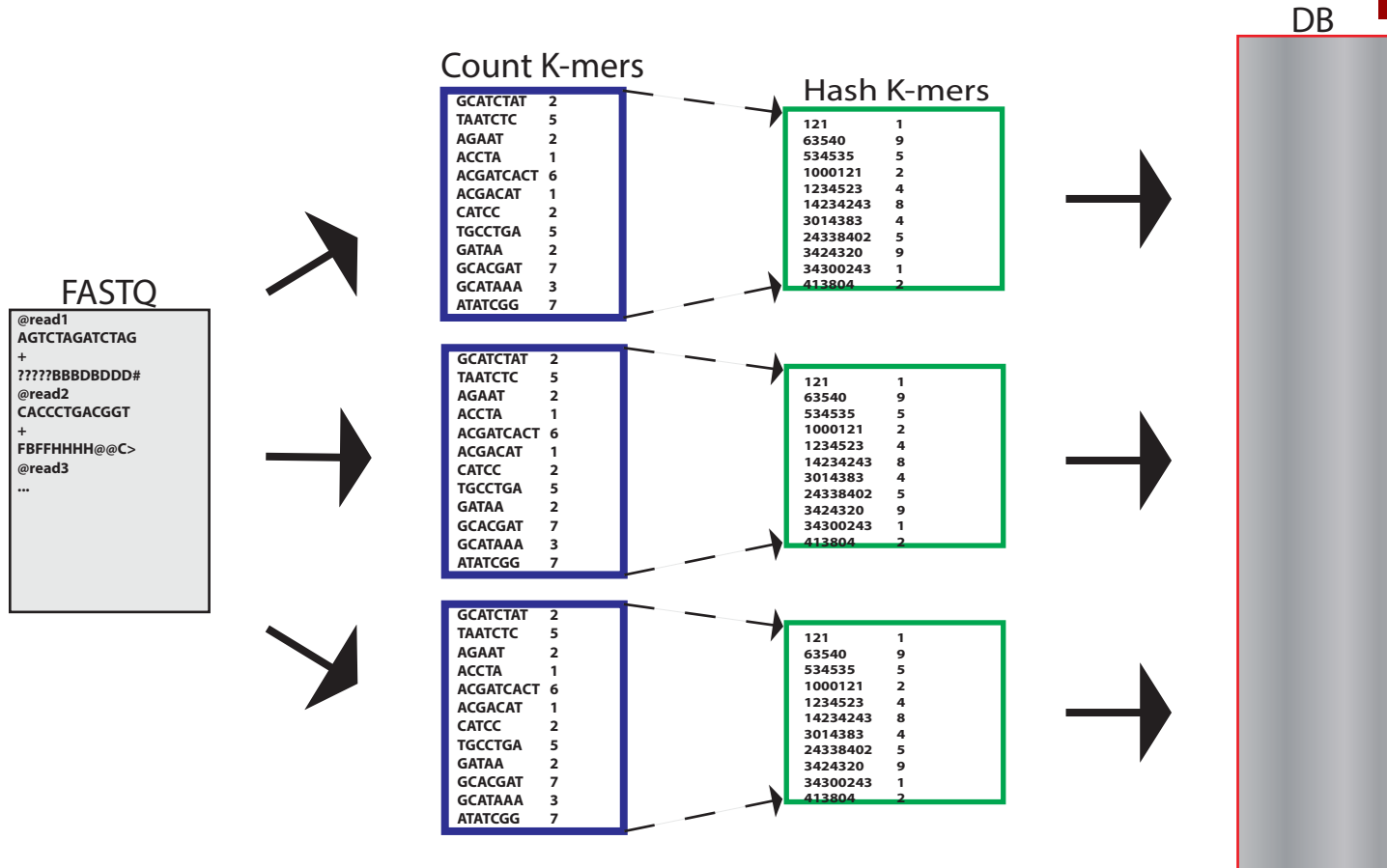
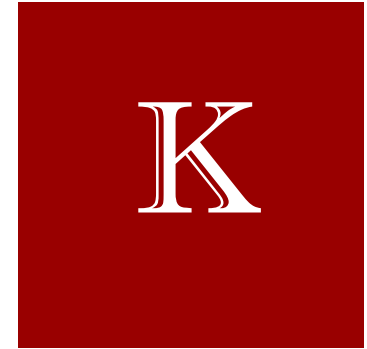


GGGTAT	17
TAATGCATC	623
AGTAA	29
ACCTAGTCA	905
ACGATCC	63
ACGTT	101
...	
CATAC	372
TGTGA	93
CTGATAG	52
GCT	37
GCATAAC	31
TTATC	706

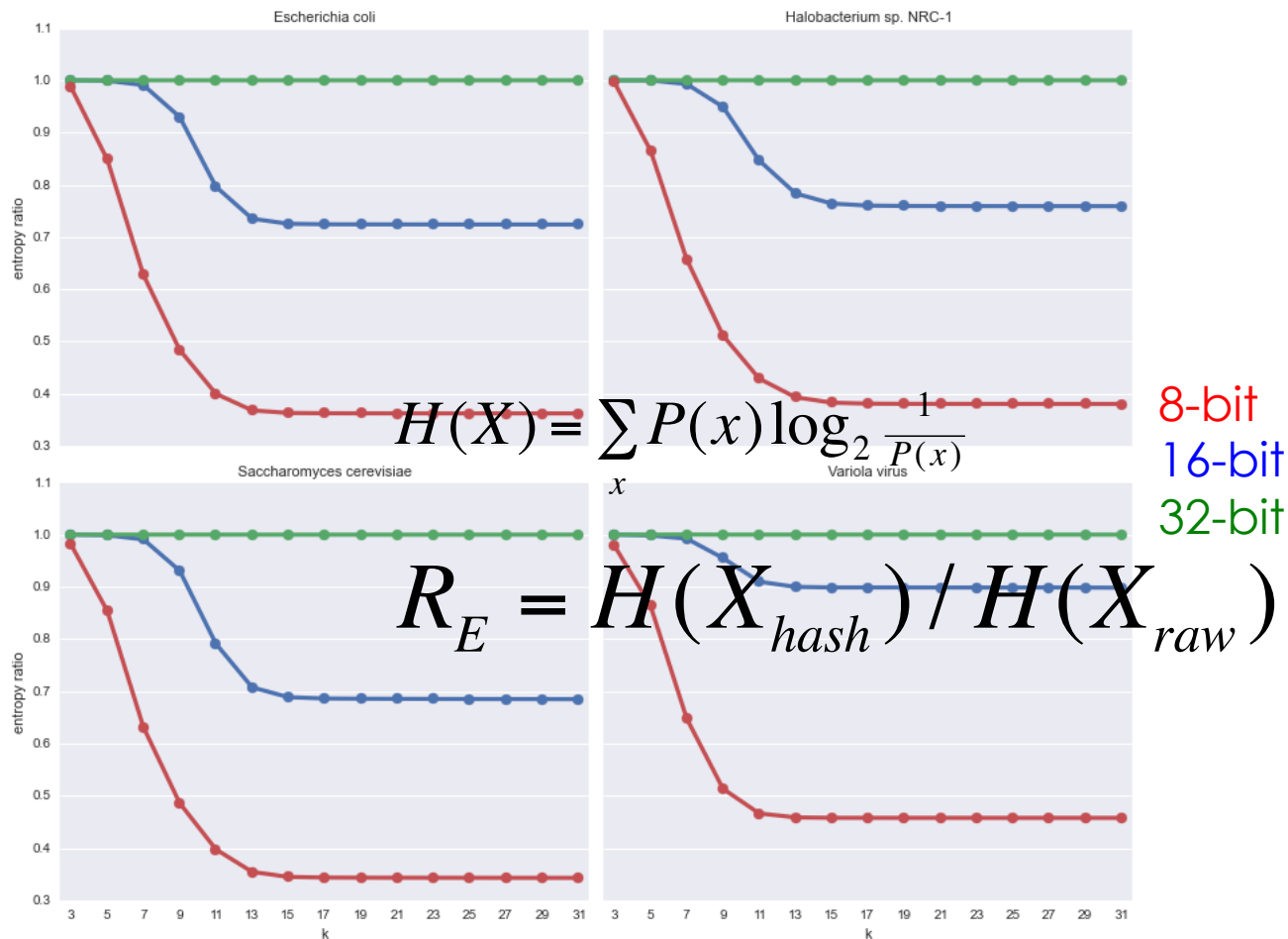
12331214	131
35350840	326
534563435	9782
30001231	1973
1232144523	31424
12234243	85
53014583	240
24348402	54545
342432	94396
4320243	1424
31124004	8937

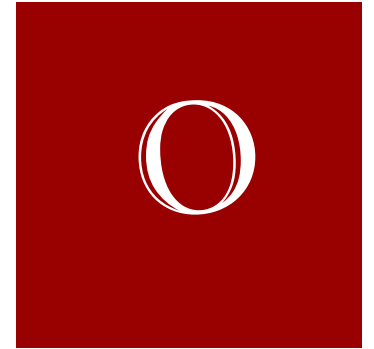


KMerge - sample



K-mer information content





Outline

- I. Introduction
- II. An alternative view of biological sequences
- III. A model for sequence identification
- IV. Conclusion/Future plans



Flipping Coins



H = Flip Coin A

T = Flip Coin B



Flipping Coins

1 H T T T H H T H T H

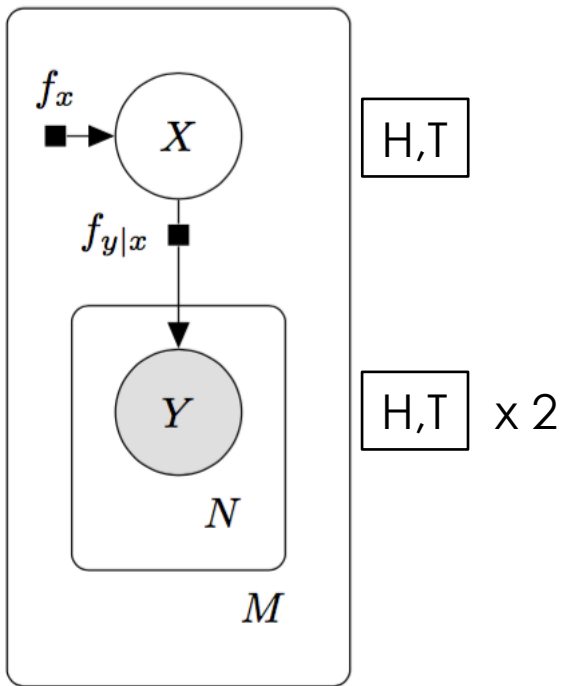
2 H H H H T H H H H H

3 H T H H H H H T H H

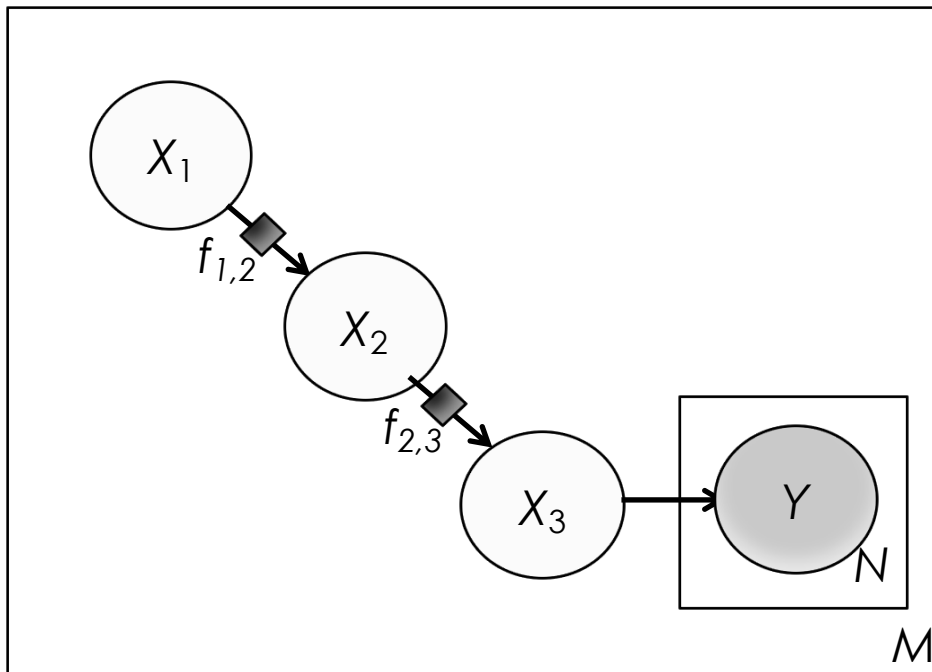
4 H T H T T T H H T T

5 T H H H T H H H T H

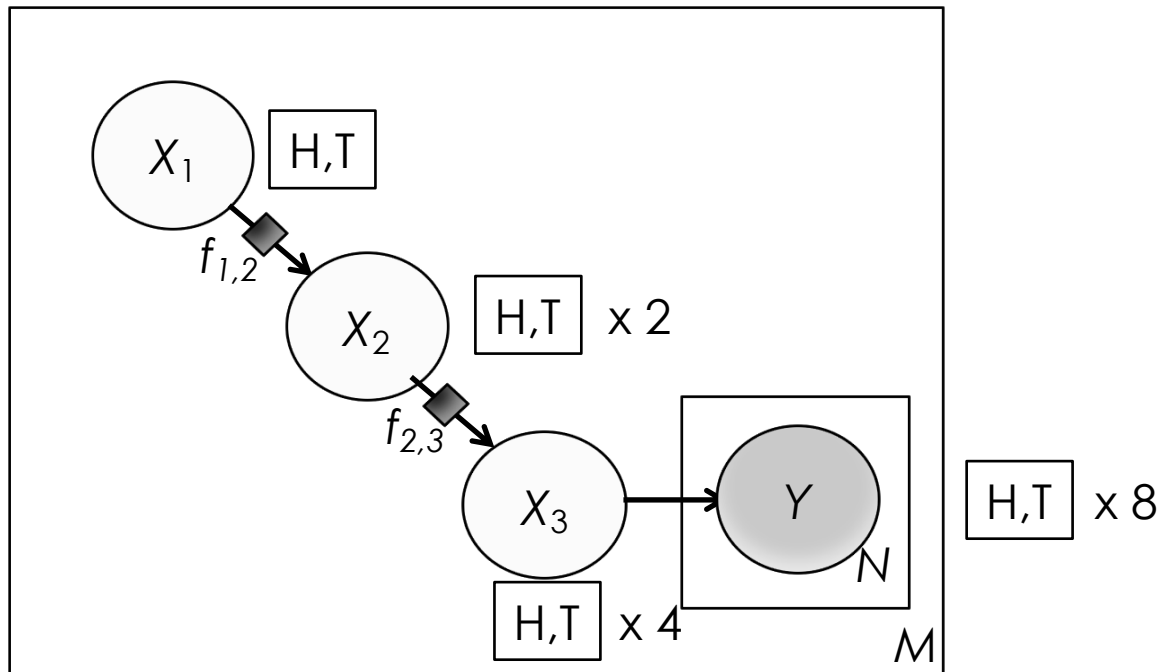
Flipping Coins



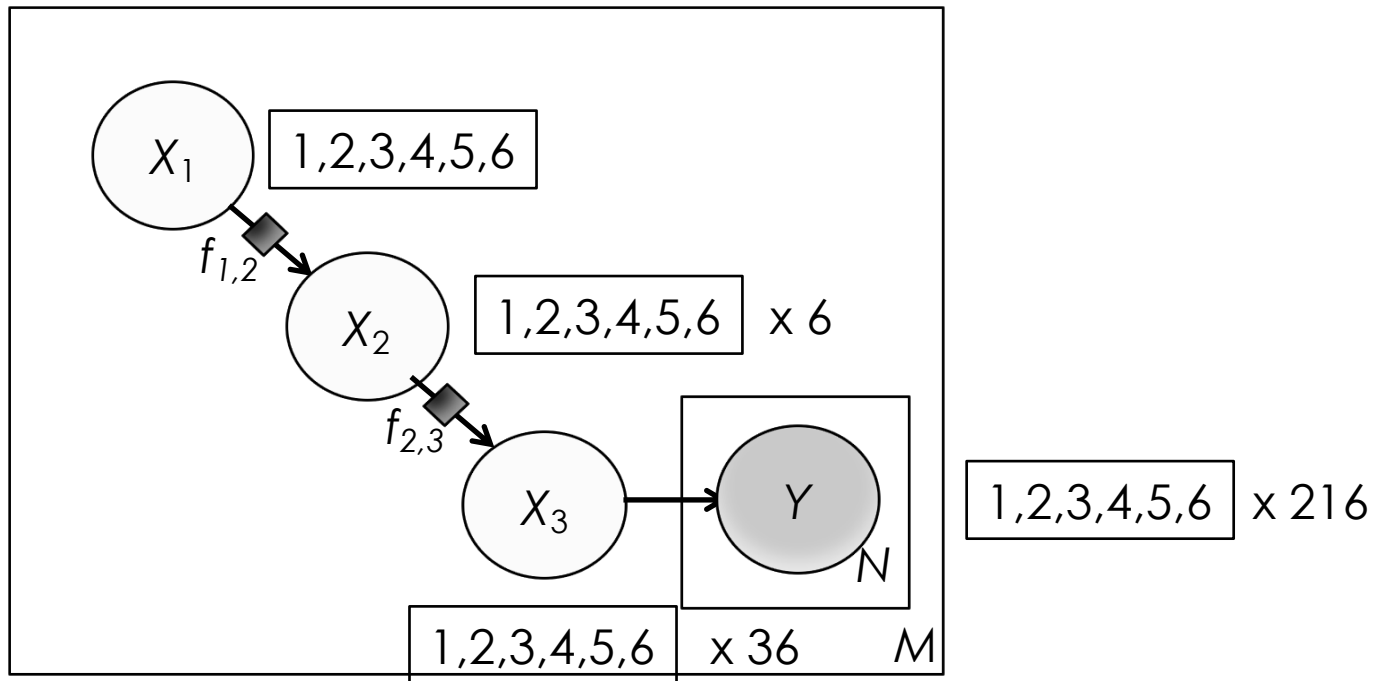
Flipping Multiple Coins



Flipping Multiple Coins



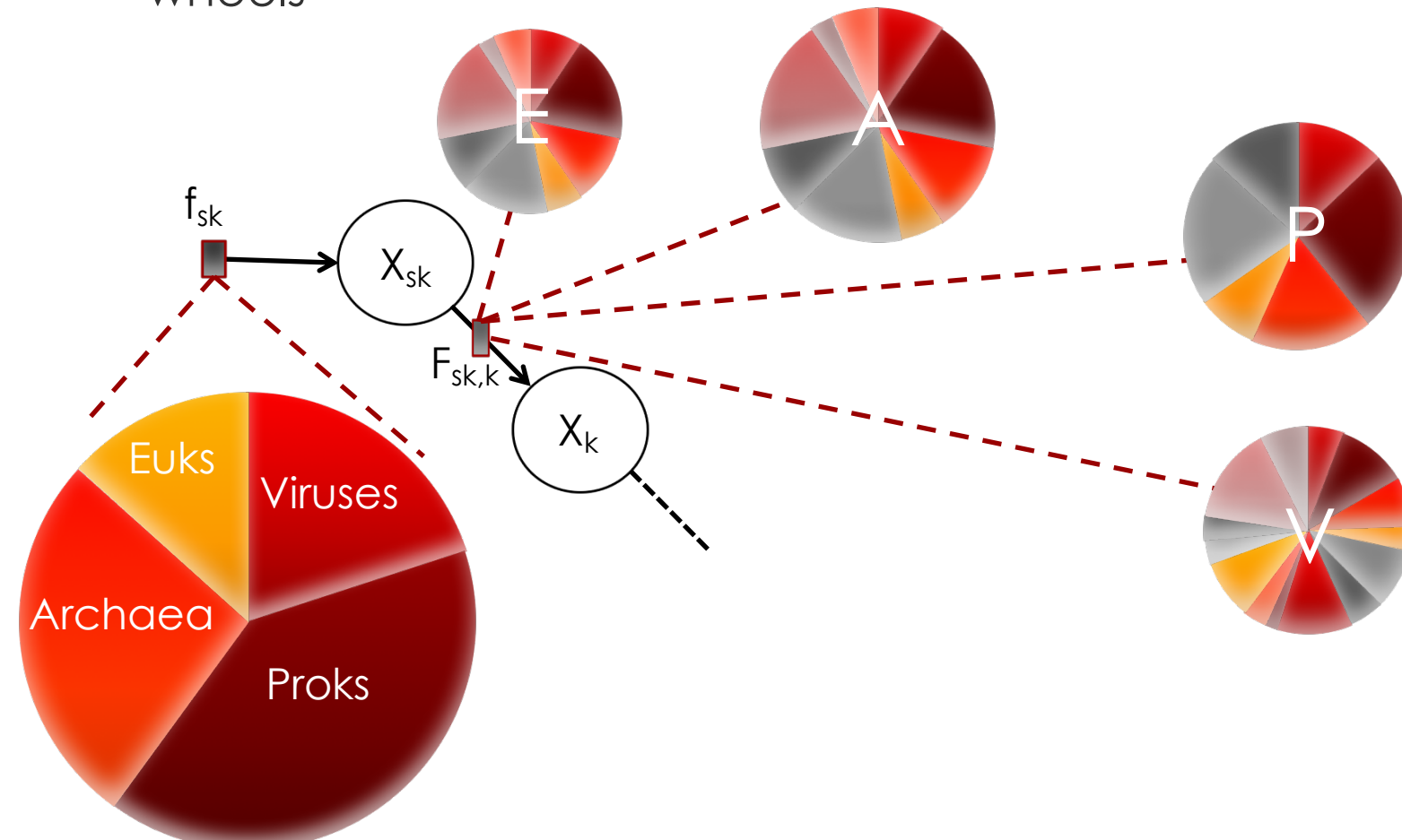
Rolling Dice



Generative model

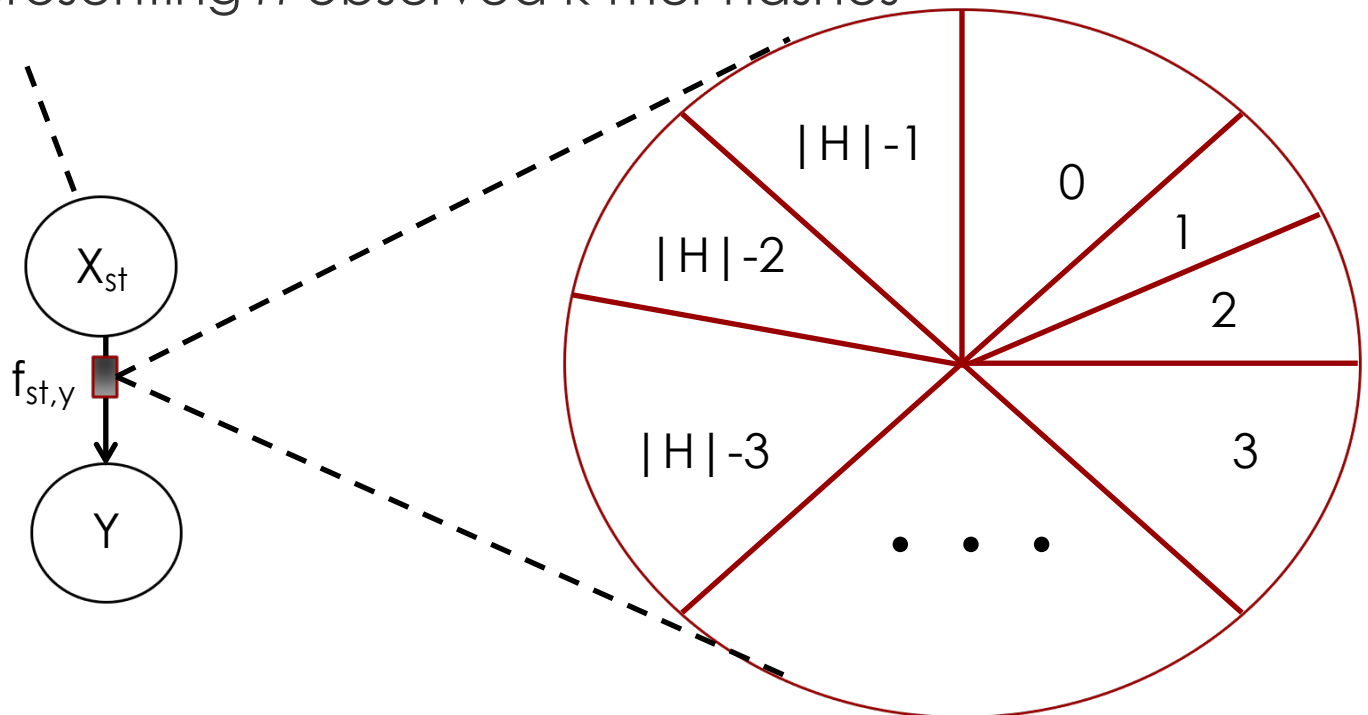


- Each factor can be envisioned as one or more sub-rank wheels

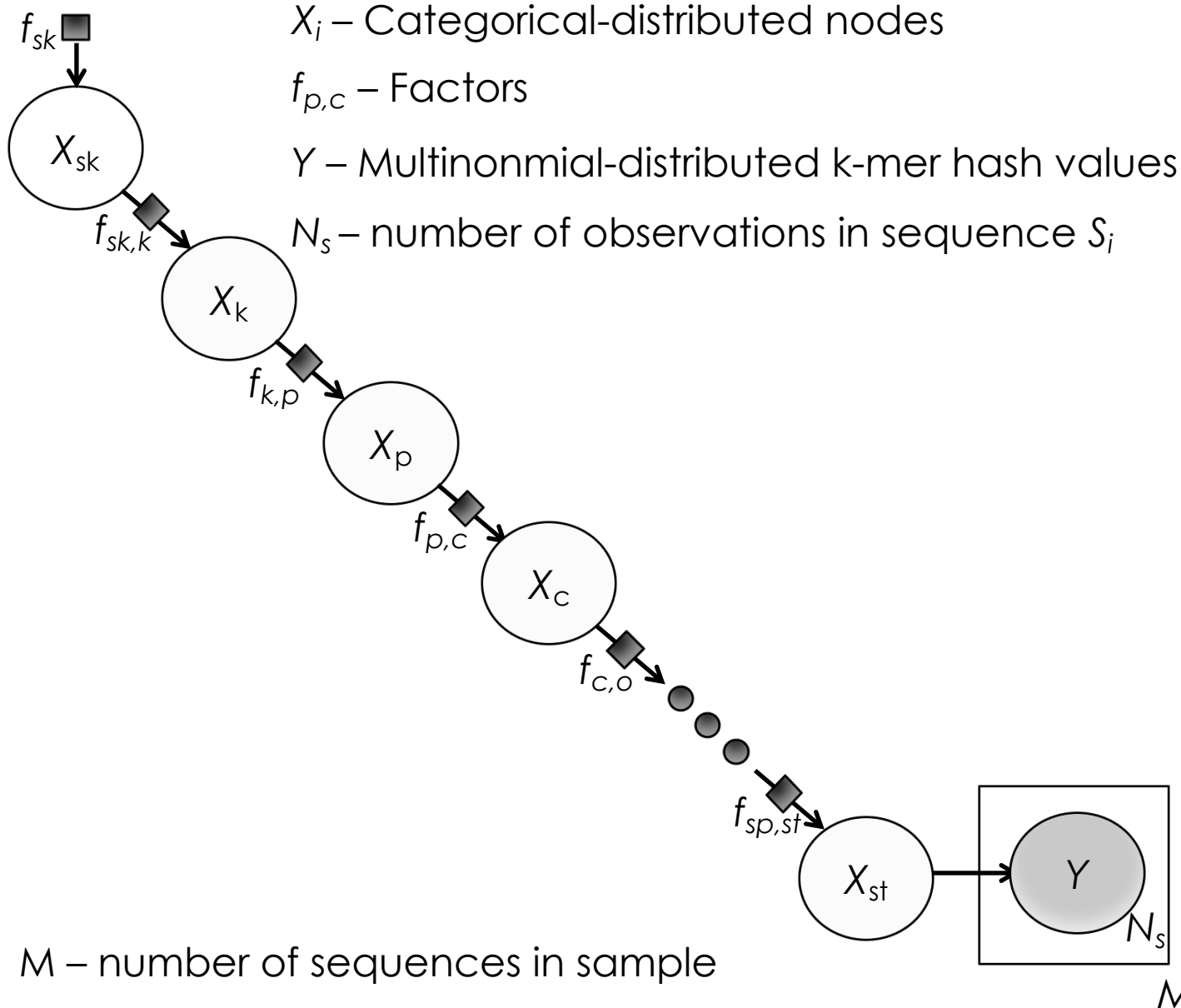


Generative process

- $|H|$ - number of possible hash values
- Hash wheel for St^* is the distribution of hashes in genome
- Hash wheel spun $N_s = n$ times
- Y_0, \dots, Y_n representing n observed k-mer hashes



Graphical model



Generative process

A red square containing a white, stylized letter 'S'.

Sequence S	
Hash Val	Count
0	4
1	0
2	9
⋮	⋮
$ H -3$	13
$ H -2$	1
$ H -1$	2



Model inference

$$Z = \{X, f\}$$

Z: hidden variables

X: rank variables

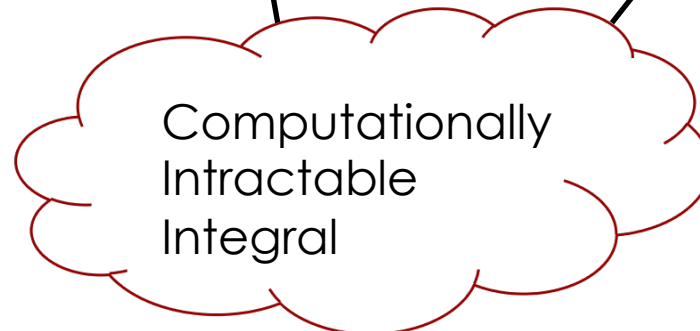
f: conditional probabilities

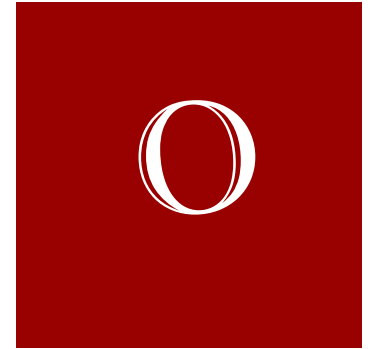
$$P(Z | Y) = ?$$

$$P(Z | Y) = \frac{P(Y | Z)P(Z)}{P(Y)}$$

Parameter Learning:

variational
approximation





Outline

- I. Introduction
- II. An alternative view of biological sequences
- III. A model for sequence identification
- IV. Conclusion/Future plans

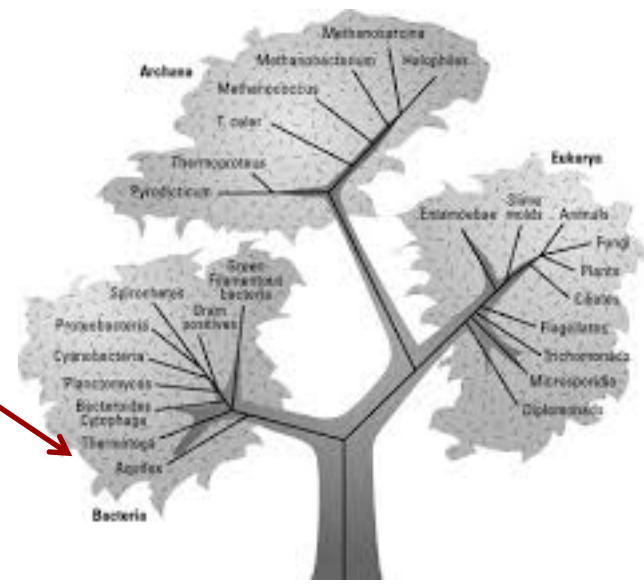
The Goal

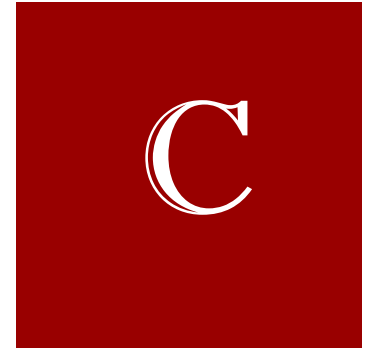


DNA Sequence



KMerge
+
Model





Next Steps

- Implement complete model
- Scaling, tuning, and benchmarking method
- Web application/visualization platform
- Publish results

Acknowledgements



Mason Lab

Chris Mason, Ph.D.*
Dhruva Chandramohan, Ph.D.*
Cem Meydan, Ph.D.*
Elizabeth Hénaff, Ph.D.*
Francine Garrett-Bakelman, MD, Ph.D.
Virginia Wagatsuma, Ph.D.
Niamh O'Hara, Ph.D.
Bharath Prithviraj, Ph.D.
Marjan Bozinoski, Ph.D.
Lenore Pipes, Ph.D.
Priya Vijay, Ph.D.
Pradeep Ambrose
Jake Reed
Alexa McIntyre
Noah Alexander
Sofia Ahsanuddin
Ebrahim Afshinnekoo
Chou Chou
Maureen Milici

Thesis Committee

Olivier Elemento, Ph.D.*
Gunnar Rättsch, Ph.D.
Jason Mezey, Ph.D.

PhD Survival Team

Noah Dukler
Tom Vincent, Ph.D.
Andre Kahles, Ph.D.
Jaakko Luttinen, Ph.D.
PBTech

Tri-I CBM

David Christini, Ph.D.
Christina Leslie, Ph.D.
Margie Hinonangan-Mendoza
Kathleen Pickering
Francine Collazo-Espinell

Funding

NIH Ruth L. Kirschstein
National Research Service
Award (F31GM111053)