

# Metagenomics workshop

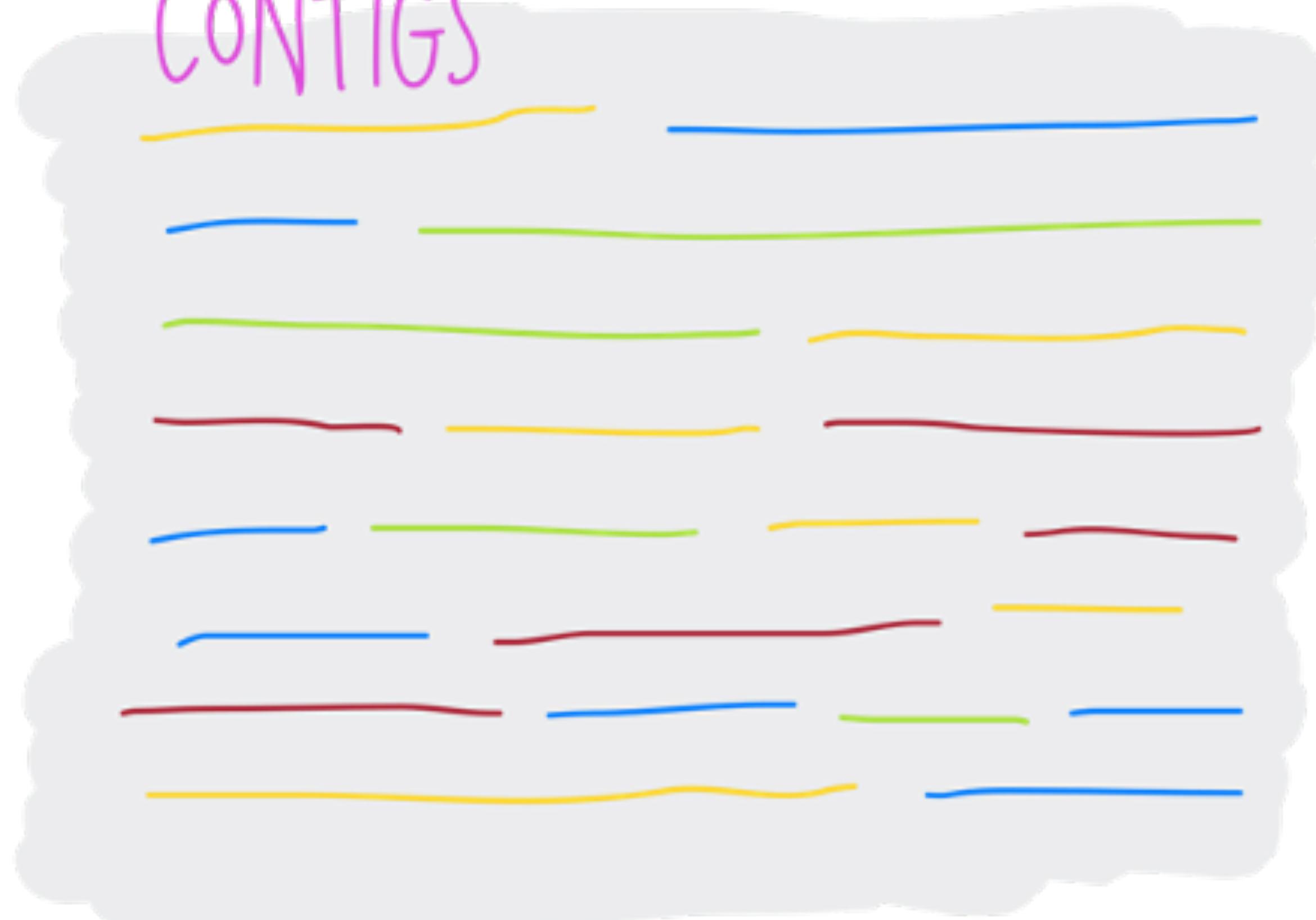
## Module 2: from bins to genomes

Lucas Paoli ([paolil@ethz.ch](mailto:paolil@ethz.ch)), Sunagawa Lab  
Institute of Microbiology, ETH Zurich

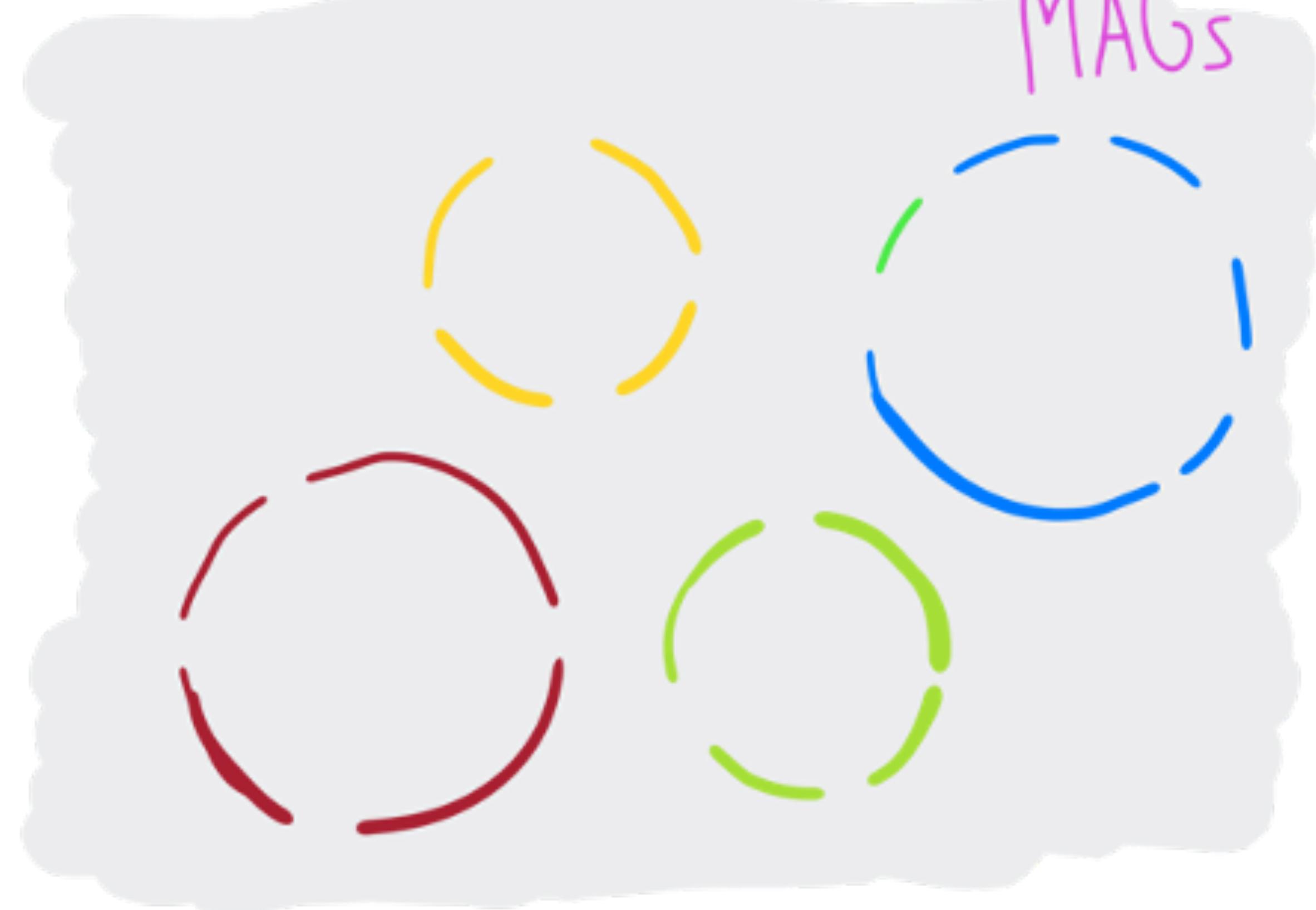
Serina Robinson ([serina.robinson@eawag.ch](mailto:serina.robinson@eawag.ch))  
Dept. of Environmental Microbiology, EAWAG

SEQUENCE COMPOSITION

CONTIGS



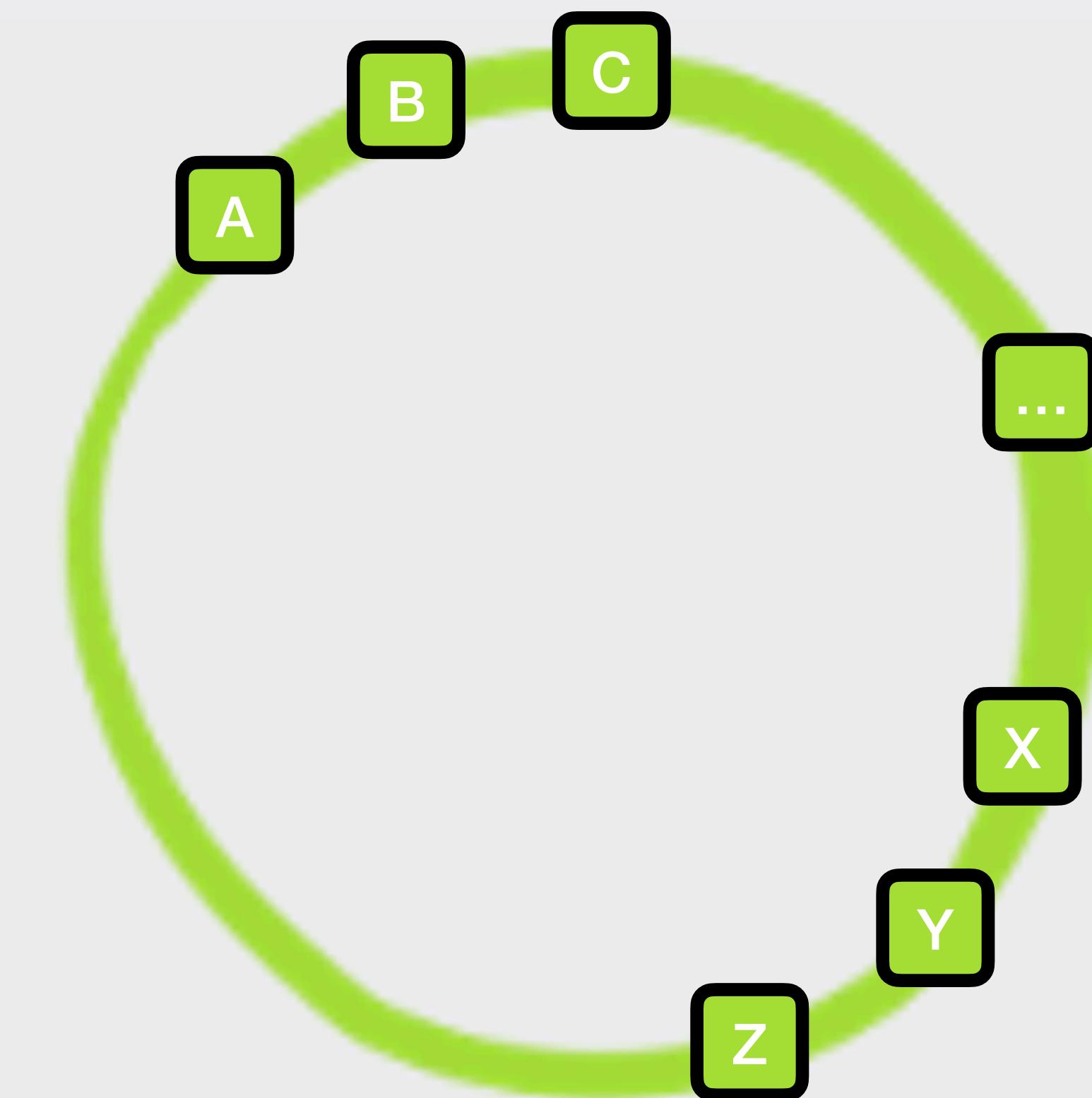
MAGs



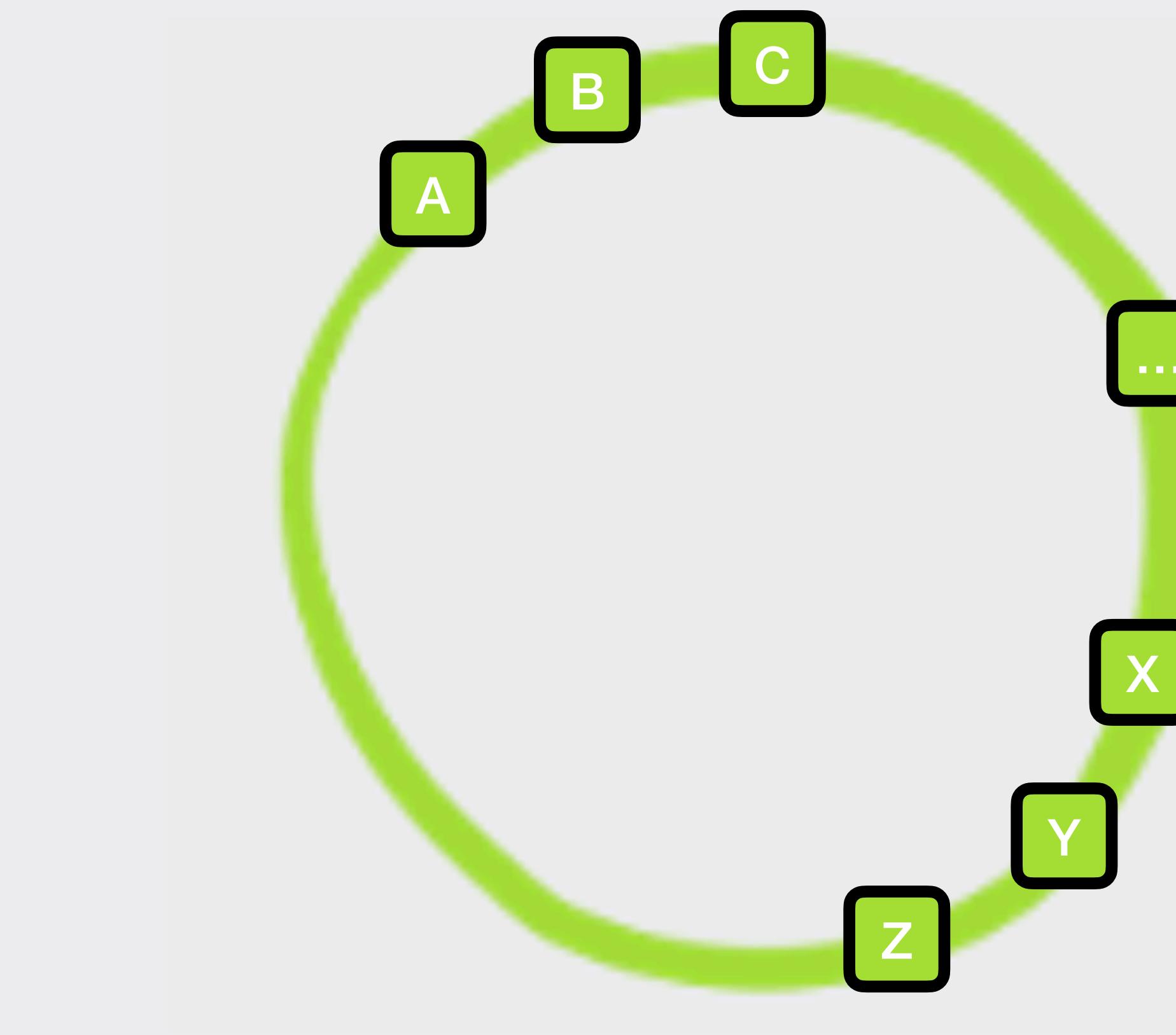
DIFFERENTIAL COVERAGE

Distinguishing spurious groups of  
contigs from draft genomes:  
completeness and contamination

# Universal single-copy marker genes

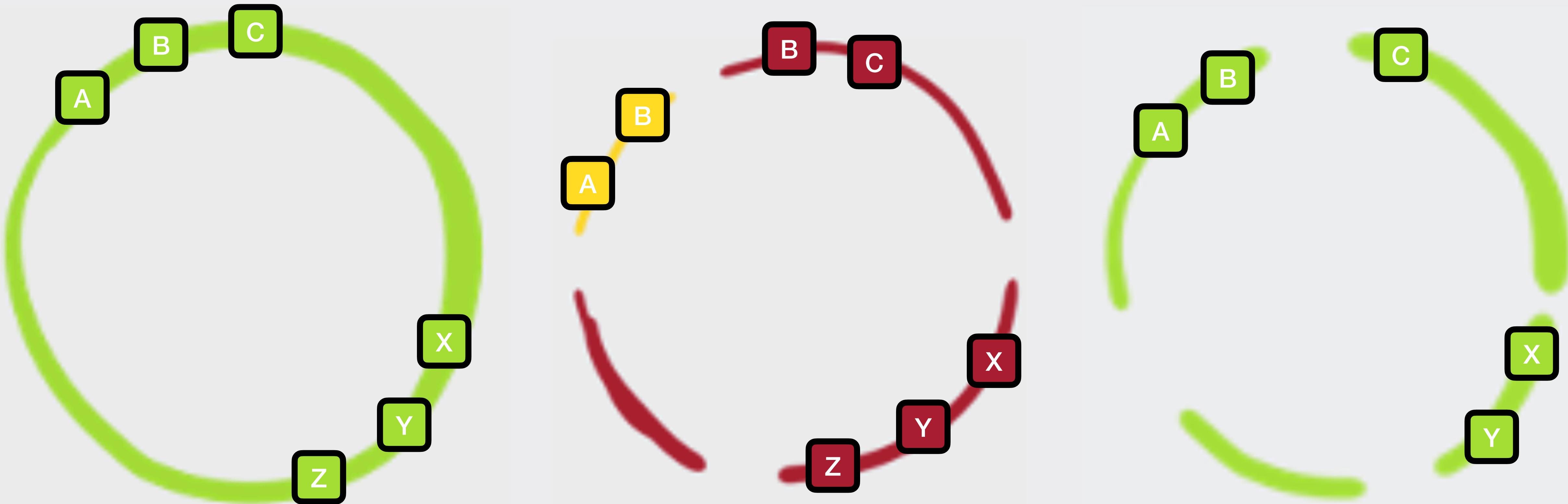


# Universal single-copy marker genes

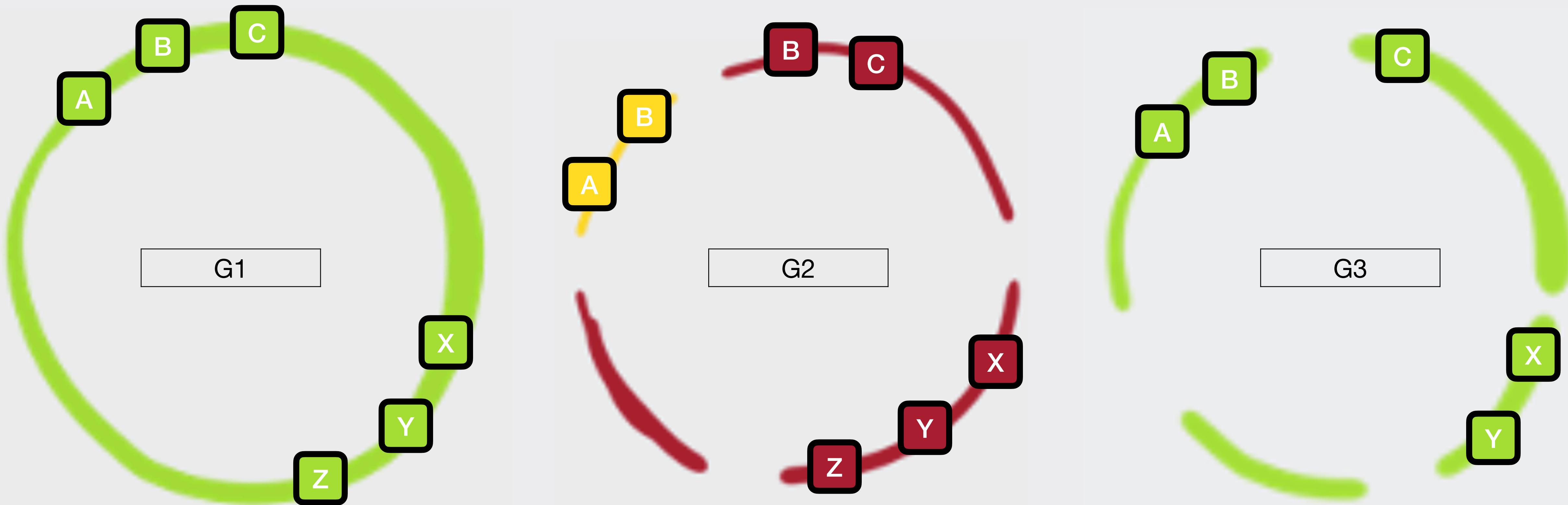


Between 40 and 120 for Bacteria/Archaea depending on cutoffs

# Universal single-copy marker genes

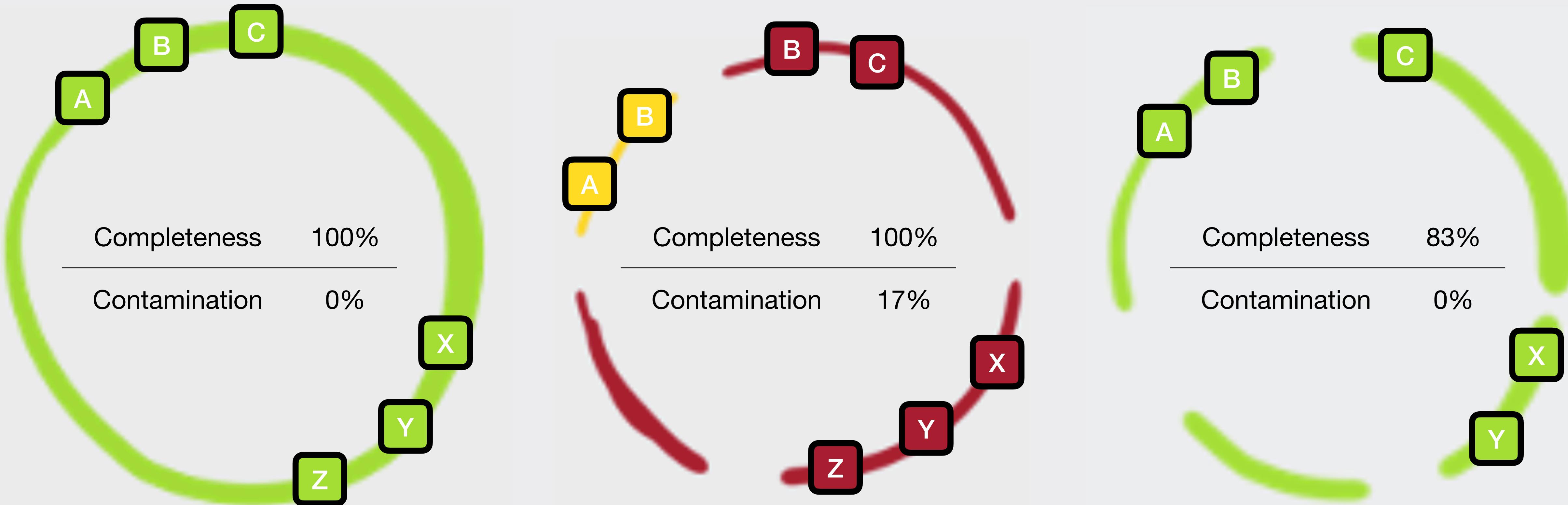


# Universal single-copy marker genes



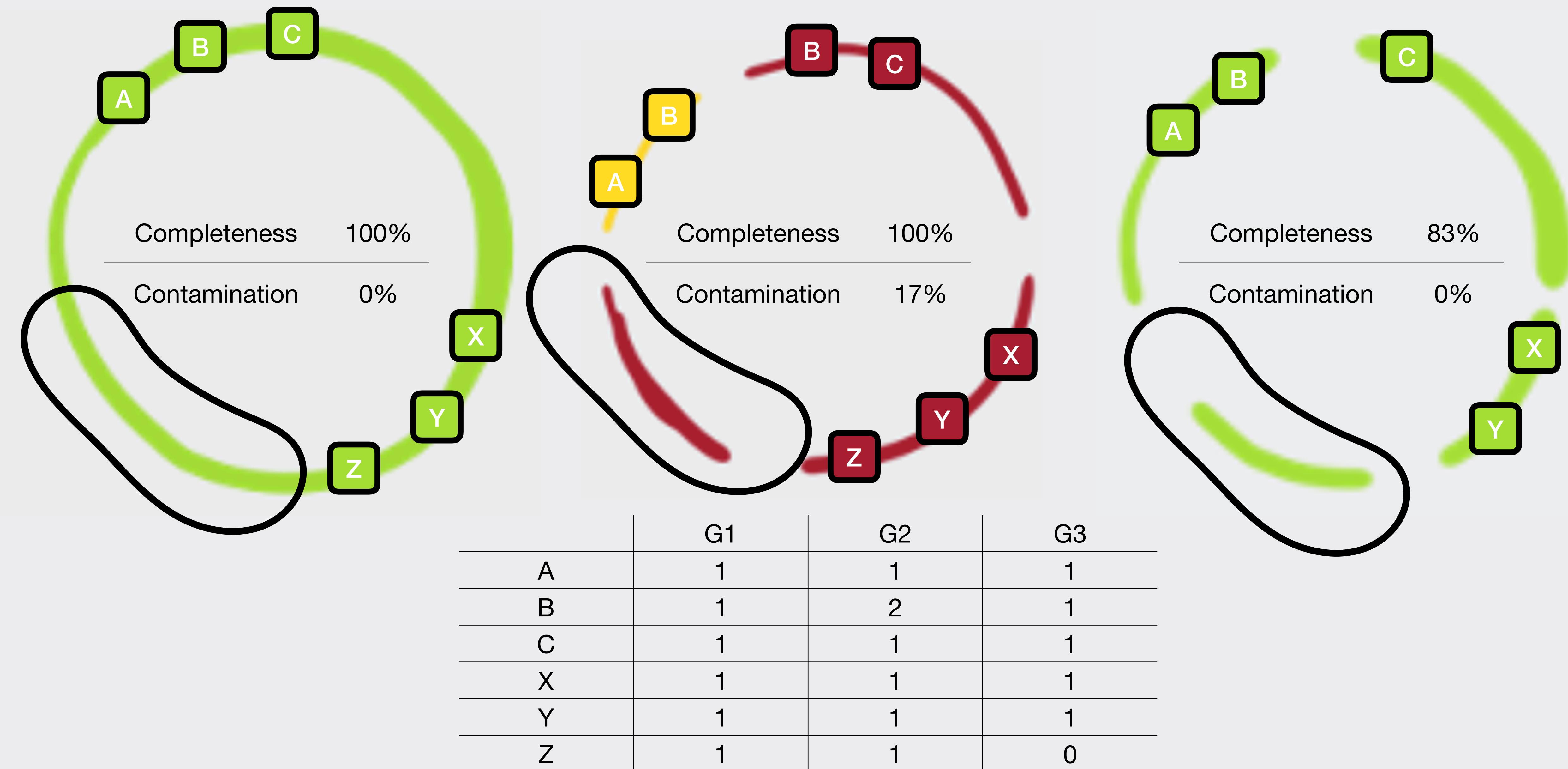
	G1	G2	G3
A	1	1	1
B	1	2	1
C	1	1	1
X	1	1	1
Y	1	1	1
Z	1	1	0

# Universal single-copy marker genes

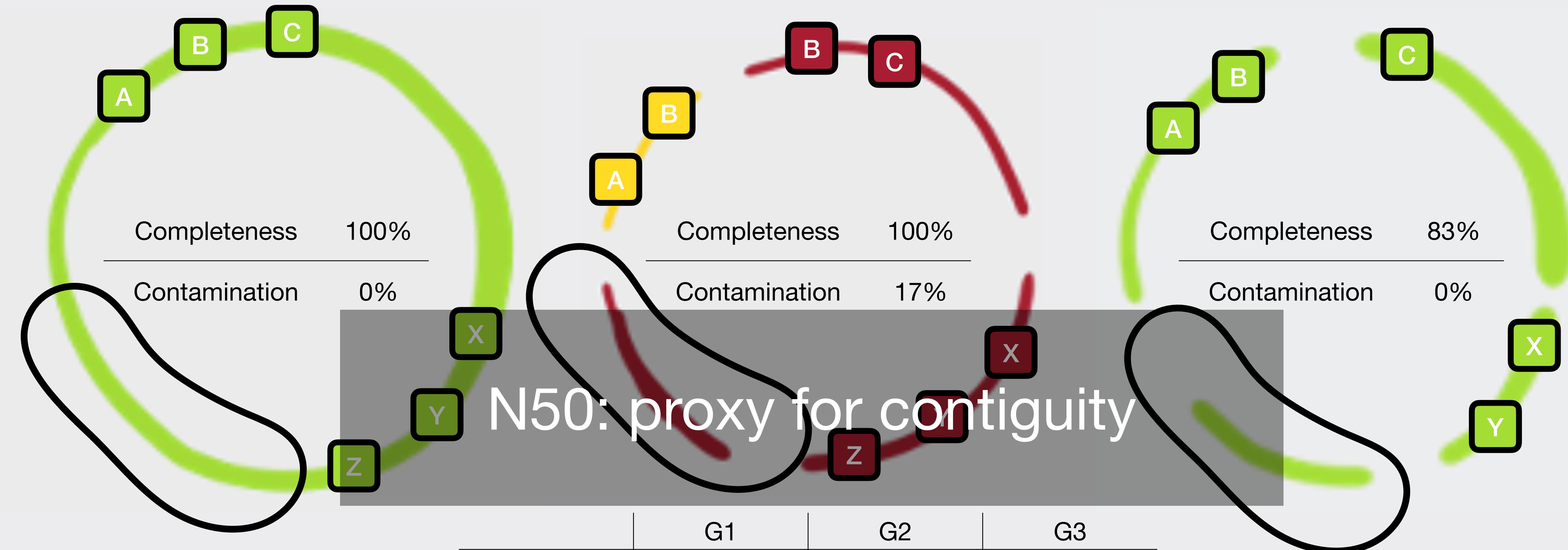


	G1	G2	G3
A	1	1	1
B	1	2	1
C	1	1	1
X	1	1	1
Y	1	1	1
Z	1	1	0

# Universal single-copy marker genes



# Universal single-copy marker genes



	G1	G2	G3
A	1	1	1
B	1	2	1
C	1	1	1
X	1	1	1
Y	1	1	1
Z	1	1	0

# In practice: evaluation software

- CheckM
- Anvi'o
- BUSCO
- ...

# Critically thinking about community standards

High-quality draft (SAG/MAG)	
Assembly quality <sup>a</sup>	Multiple fragments where gaps span repetitive regions. Presence of the 23S, 16S, and 5S rRNA genes and at least 18 tRNAs.
Completion <sup>b</sup>	>90%
Contamination <sup>c</sup>	<5%
Medium-quality draft (SAG/MAG)	
Assembly quality <sup>a</sup>	Many fragments with little to no review of assembly other than reporting of standard assembly statistics.
Completion <sup>b</sup>	≥50%
Contamination <sup>c</sup>	<10%

nature  
biotechnology  
OPEN

PERSPECTIVE

## Minimum information about a single amplified genome (MISAG) and a metagenome-assembled genome (MIMAG) of bacteria and archaea

Robert M Bowers<sup>1</sup>, Nikos C Kyrpides<sup>1</sup>, Ramunas Stepanauskas<sup>2</sup> , Miranda Harmon-Smith<sup>1</sup>, Devin Doud<sup>1</sup>, T B K Reddy<sup>1</sup>, Frederik Schulz<sup>1</sup> , Jessica Jarett<sup>1</sup>, Adam R Rivers<sup>1,3</sup>, Emiley A Eloe-Fadrosh<sup>1</sup>, Susannah G Tringe<sup>1,4</sup> , Natalia N Ivanova<sup>1</sup>, Alex Copeland<sup>1</sup>, Alicia Clum<sup>1</sup>, Eric D Becraft<sup>2</sup>, Rex R Malmstrom<sup>1</sup>, Bruce Birren<sup>5</sup>, Mircea Podar<sup>6</sup>, Peer Bork<sup>7</sup>, George M Weinstock<sup>8</sup>, George M Garrity<sup>9</sup>, Jeremy A Dodsworth<sup>10</sup>, Shibu Yooseph<sup>11</sup>, Granger Sutton<sup>12</sup> , Frank O Glöckner<sup>13</sup>, Jack A Gilbert<sup>14,15</sup>, William C Nelson<sup>16</sup>, Steven J Hallam<sup>17</sup>, Sean P Jungbluth<sup>1,18</sup> , Thijs J G Ettema<sup>19</sup>, Scott Tighe<sup>20</sup>, Konstantinos T Konstantinidis<sup>21</sup>, Wen-Tso Liu<sup>22</sup>, Brett J Baker<sup>23</sup>, Thomas Rattei<sup>24</sup>, Jonathan A Eisen<sup>25</sup>, Brian Hedlund<sup>26,27</sup>, Katherine D McMahon<sup>28,29</sup>, Noah Fierer<sup>30,31</sup>, Rob Knight<sup>32</sup> , Rob Finn<sup>33</sup>, Guy Cochrane<sup>33</sup>, Ilene Karsch-Mizrachi<sup>34</sup>, Gene W Tyson<sup>35</sup>, Christian Rinke<sup>35</sup> , The Genome Standards Consortium<sup>36</sup>, Alla Lapidus<sup>37</sup> , Folker Meyer<sup>14</sup>, Pelin Yilmaz<sup>13</sup> , Donovan H Parks<sup>35</sup> , A Murat Eren<sup>38</sup> , Lynn Schriml<sup>39</sup>, Jillian F Banfield<sup>40</sup>, Philip Hugenholtz<sup>35</sup> & Tanja Woyke<sup>1,4</sup>

The background of the slide features a dark, abstract design. A glowing, translucent DNA double helix is centered, with its bright white and yellow segments standing out against the dark blue and black background. Several small, glowing liquid droplets are scattered across the surface, adding to the organic feel. The overall aesthetic is scientific and modern.

Taxonomic annotation

# Genome-based classification and taxonomy

## RESOURCE

nature  
biotechnology

A standardized bacterial taxonomy based on genome phylogeny substantially revises the tree of life

Donovan H Parks, Maria Chuvochina, David W Waite, Christian Rinke<sup>ID</sup>, Adam Skarszewski, Pierre-Alain Chaumeil & Philip Hugenholtz<sup>ID</sup>

Taxonomy is an organizing principle of biology and is ideally based on evolutionary relationships among organisms. Development of a robust bacterial taxonomy has been hindered by an inability to obtain most bacteria in pure culture and, to a lesser extent, by the historical use of phenotypes to guide classification. Culture-independent sequencing technologies have matured sufficiently that a comprehensive genome-based taxonomy is now possible. We used a concatenated protein phylogeny as the basis for a bacterial taxonomy that conservatively removes polyphyletic groups and normalizes taxonomic ranks on the basis of relative evolutionary divergence. Under this approach, 58% of the 94,759 genomes comprising the Genome Taxonomy Database had changes to their existing taxonomy. This result includes the description of 99 phyla, including six major monophyletic units from the subdivision of the Proteobacteria, and amalgamation of the Candidate Phyla Radiation into a single phylum. Our taxonomy should enable improved classification of uncultured bacteria and provide a sound basis for ecological and evolutionary studies.

Bioinformatics, 36(6), 2020, 1925–1927

doi: 10.1093/bioinformatics/btz848

Advance Access Publication Date: 15 November 2019

Applications Note

OXFORD

Genome analysis

**GTDB-Tk: a toolkit to classify genomes with the Genome Taxonomy Database**

Pierre-Alain Chaumeil\*, Aaron J. Mussig<sup>ID</sup>, Philip Hugenholtz and Donovan H. Parks\*

Australian Centre for Ecogenomics, School of Chemistry and Molecular Biosciences, The University of Queensland, St Lucia, QLD 4072, Australia

\*To whom correspondence should be addressed.

Associate Editor: John Hancock

Received on July 12, 2019; revised on October 15, 2019; editorial decision on November 11, 2019; accepted on November 13, 2019

## RESOURCE

<https://doi.org/10.1038/s41564-022-01214-9>

nature  
microbiology

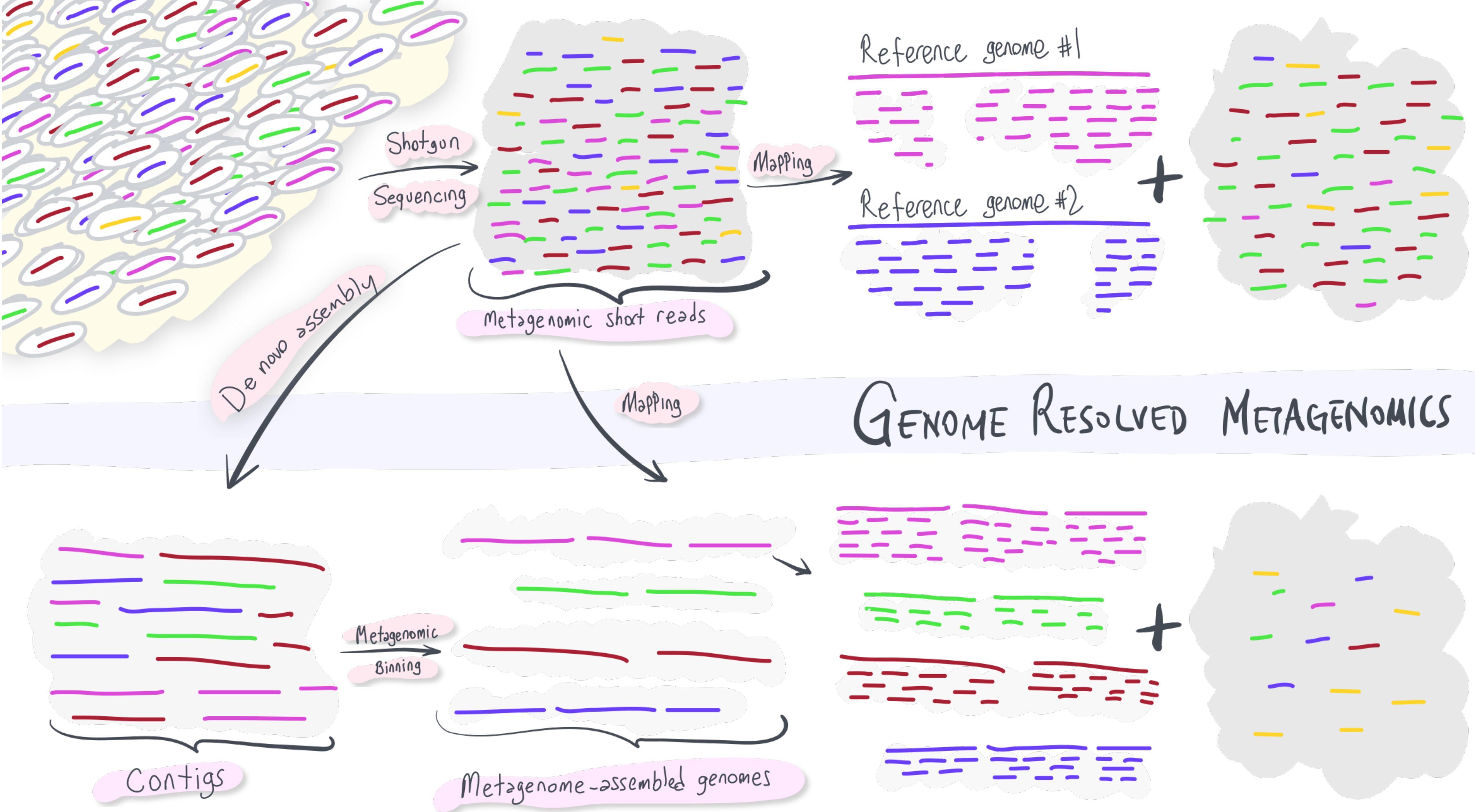


**OPEN**  
**SeqCode: a nomenclatural code for prokaryotes described from sequence data**

Brian P. Hedlund<sup>ID</sup><sup>1</sup>, Maria Chuvochina<sup>2</sup>, Philip Hugenholtz<sup>2</sup>, Konstantinos T. Konstantinidis<sup>ID</sup><sup>3</sup>, Alison E. Murray<sup>ID</sup><sup>4</sup>, Marike Palmer<sup>ID</sup><sup>1</sup>, Donovan H. Parks<sup>ID</sup><sup>2</sup>, Alexander J. Probst<sup>5</sup>, Anna-Louise Reysenbach<sup>ID</sup><sup>6</sup>, Luis M. Rodriguez-R<sup>ID</sup><sup>7</sup>, Ramon Rossello-Mora<sup>ID</sup><sup>8</sup>, Iain C. Sutcliffe<sup>ID</sup><sup>9</sup>, Stephanus N. Venter<sup>ID</sup><sup>10</sup> and William B. Whitman<sup>ID</sup><sup>11</sup>✉

Most prokaryotes are not available as pure cultures and therefore ineligible for naming under the rules and recommendations of the International Code of Nomenclature of Prokaryotes (ICNP). Here we summarize the development of the SeqCode, a code of nomenclature under which genome sequences serve as nomenclatural types. This code enables valid publication of names of prokaryotes based upon isolate genome, metagenome-assembled genome or single-amplified genome sequences. Otherwise, it is similar to the ICNP with regard to the formation of names and rules of priority. It operates through the SeqCode Registry (<https://seqco.de/>), a registration portal through which names and nomenclatural types are registered, validated and linked to metadata. We describe the two paths currently available within SeqCode to register and validate names, including *Candidatus* names, and provide examples for both. Recommendations on minimal standards for DNA sequences are provided. Thus, the SeqCode provides a reproducible and objective framework for the nomenclature of all prokaryotes regardless of cultivability and facilitates communication across microbiological disciplines.

# Cultivation-independent genome-resolved metagenomics: a summary



# Material

This course uses a lot of material from <https://merenlab.org/momics/>, I invite you to have a look.

If you want details on the bioinformatics behind you can start by having a look here:  
[https://astrobiomike.github.io/genomics/metagen\\_anvio](https://astrobiomike.github.io/genomics/metagen_anvio)



**A. Murat Eren (Meren) (PI)**

- Web  Email  Twitter  LinkedIn  Github  ORCiD
- Address: Knapp Center for Biomedical Discovery, 900 E. 57th St., MB 9, RM 9118, Chicago, IL 60637 USA
- Phone: +1-773-702-5935  Fax: +1-773-702-2281

*I am a computer scientist with a deep appreciation for the complexity of life. I design algorithms and experiments to better understand microbes and their ecology. [photos: 1, 2, 3].*

- » MBL Fellow, [Marine Biological Laboratory](#).
- » Assistant Professor, [The Department of Medicine at the University of Chicago](#).
- » Committee on Microbiology, [The Biomedical Sciences Cluster at the University of Chicago](#).



**Mike Lee**

- Web  Email  Twitter  LinkedIn  Github
- » NASA Space Biology Fellow, [NASA Ames Research Center](#).
- » JCVI Research Fellow, [J. Craig Venter Institute](#).

- 👉 Combining reference genome annotations with your own in pangenomes (Sat, Dec 01, 2018)
- 👉 Anvi'o 'views' demystified (Mon, May 08, 2017)
- 👉 Making anvi'o use your own HMM collection (Sat, May 21, 2016)



# Questions?

Image: François Aurat