

Introduction to shotgun meta'omic analysis

Eric A. Franzosa, Ph.D. Galeb Abu-Ali, Ph.D.

Harvard CFAR Workshop on Metagenomics 17 September 2015



Huttenhower Research Group Harvard Chan School of Public Health Department of Biostatistics





- Shotgun meta'omics primer
- Informal survey
- Meta'omic taxonomic profiling
 - MetaPhlAn(2)
- Meta'omic functional profiling
 - Broad functional profiling with HUMAnN(2)
 - Targeted functional profiling with ShortBRED
 - Predictive functional profiling with PICRUSt
- Downstream analyses
- Resources
- Tutorials (today and tomorrow)

Sequencing as a tool for microbial community analysis (amplicon vs. shotgun)



A note on "metagenomics" vocabulary

Amplicon sequencing

DOC

- PCR amplify and seq. specific marker(s)
- Often the 16S rRNA gene (for bacteria)
- Shotgun sequencing
 - Seq. short, random DNA/RNA fragments
 - Whole metagenome shotgun (WMS)
 - Whole metatranscriptome shotgun
 - Collectively, meta'omic sequencing

Sequencing as a tool for microbial community analysis (amplicon vs. shotgun)

- Where they overlap
 - Strengths
 - Quantifying taxonomic abundance
 - Ecological statistics
 - Taxon-taxon association
 - Taxon-metadata association
 - Challenges
 - Compositional (& noisy) data
 - Difficult distributions
 - Biases from sequencing

Sequencing as a tool for microbial community analysis (amplicon vs. shotgun)

- Properties of shotgun meta'omic sequencing
 - Strengths
 - Taxonomic resolution (species, strains)
 - Functional genomics (genes, transcripts)
 - Comparative genomics
 - Challenges
 - More expensive per sample
 - Data are bigger, compute more intensive
 - Need a good reference
 - Contamination

Survey: who has/plans to work with...

• 16S data?

DOC

- shotgun data?
- metatranscriptomes?
- human vs. environmental samples?
- metagenomic assemblies?

The universal meta'omics workflow



We develop computational methods in these areas Today we'll be focusing on this subset

Meta'omic quality control

- Trim low-quality bases from read ends
 - http://www.usadellab.org/cms/?page=trimmomatic
- Drop short reads

DOC

- Remove contaminant sequences
 - E.g. human genome, EST database
- Remove low-complexity sequences (?)
- Enforce end-pairing (?)
 - Not required for bioBakery tools
- Integrated workflow coming soon!
 - https://bitbucket.org/biobakery/kneaddata



Who is there? (taxonomic profiling)

What are they doing? (functional profiling)

Slide by Dirk Gevers

The NIH Human Microbiome Project (HMP): A comprehensive microbial survey

Nasal

Oral

Skin

Gastrointestinal

Urogenital

- What is a "normal" human microbiome?
- 300 healthy human subjects
- Multiple body sites
 - 15 male, 18 female
- Multiple visits
- Clinical metadata

www.hmpdacc.org



Profiling microbial communities and ecology at species-level resolution (HMP)



12



tongue dorsum

attached keratinized gingiya

stoo

supragingival plague

pauricular crease

vaginal introitus

mid vagina

Are there discrete "types" of typical human microbiomes?



14



MetaPhlAn(2) For meta'omic taxonomic profiling



RepoPhlAn

ChocoPhIAn (http://metaref.org)









Evaluation of MetaPhlAn accuracy

DOC



(Validation on high-complexity uniformly distributed synthetic metagenomes.)







- MetaPhlAn 1.0 focused on bacteria and archaea
- v2.0 adds support for eukaryotes and viruses
- In along with many more bacteria and archaea
- v2.0 supports profiling at the strain level

DOC MetaPhIAn2: synthetic evaluation



https://bitbucket.org/biobakery/metaphlan2 23

MetaPhlAn2: Results for HMP Skin

DOC



MetaPhIAn in action: strain profiling



- In practice, not all markers are present
- Individual-specific marker "barcodes"
- Often very stable over time

DOC



Who is there? (taxonomic profiling)

What are they doing? (functional profiling)

(What we mean by "function")

INOSITOL PHOSPHATE METABOLISM



00562 11/1/10 (c) Kanehisa Laboratories

Metagenomic analyses: gene calling and proxygenes





Orphelia: Hoff, 2009 MetaGene: Noguchi, 2006

DOC





Orthology: Grouping genes by conserved sequence features COG, KO, FIGfam...



Structure: Grouping genes by similar protein domains Pfam, TIGRfam, SMART, EC...

Biological roles:

Grouping genes by pathway and process involvement GO, KEGG, MetaCyc, SEED...



Warnecke, 2007



<u>DeLong</u>, 2006

в

Niche specialization in human microbiome function

LEfSe: <u>LDA Effect Size</u> Nonparametric test for microbial and metagenomic biomarkers <u>http://huttenhower.sph.harvard.edu/lefse</u>



Most processes are "core": <10% are differentially present/absent even by body site
 Contrast zero microbes meeting this threshold!

• Most *processes are habitat-adapted*: >66% are *differentially abundant* by body site

Um



"Who's there," versus, "What they're doing," in the healthy human microbiome $_{\leftarrow Subjects \rightarrow}$



http://hmpdacc.org/HMMRC

Which *functions* of the gut microbiome are disrupted by IBD?

- Over <u>six times</u> as many microbial metabolic processes disrupted in IBD as microbes.
 - If there's a transit strike, everyone working for the MBTA is disrupted, not everyone named Smith or Jones.
 - Phylogenetic distribution of function is *consistent* but *diffuse*
- During IBD, microbes...

Stop

- Creating most amino acids
- Degrading complex carbs.

DOC

Producing short-chain fatty acids

Start

- Taking up more host products
- Dodging the immune system
- Adhering to and invading host cells

Integrated functional meta'omics (Examining community DNA & RNA)



Functional metagenomics & metatranscriptomics of 8 heathy human stool samples

Franzosa et al. PNAS 11:E2329-38 (2014)



HUMAnN(2)

For broad meta'omic functional profiling



<u>HMP Unified Metabolic Analysis Network</u>



HUMAnN

DOC

Sample 1Sample 2Sample 3Sample 4Sample 5AImage: Sample 2Image: Sample 3Image: Sample 4Image: Sample 5BImage: Sample 3Image: Sample 3Image: Sample 4Image: Sample 5CImage: Sample 3Image: Sample 3Image: Sample 3Image: Sample 4CImage: Sample 3Image: Sample 3Image: Sample 3Image: Sample 3

Short reads + protein families Translated BLAST search

$$c(g) = \frac{1}{|g|} \sum_{r} \frac{\sum_{a(r)} (1 - p_a) \Delta(a = g)}{\sum_{a(r)} (1 - p_a)}$$

Weight hits by significance

Sum over families

Adjust for sequence length

Repeat for each metagenomic or metatranscriptomic sample





Millions of hits are collapsed into thousands of gene families (*still a large number*)



Map genes to pathways

- Use MinPath (Ye 2009) to find simplest pathway explanation for observed genes
- Remove pathways unlikely to be present due to low organismal abundance
- Smooth/fill gaps

Collapsing gene family abundance into pathway abundance (or presence/absence) yields a smaller, more tractable feature set





Validated against synthetic metagenome samples (similar to MetaPhlAn validation)

Gene family abundance and pathway presence/absence calls beat naïve best-BLAST-hit strategy



- Avoid translated search where possible
- Speed up translated search with ORF-picking
- Stratify community-wide function by organism
- Focus on open gene family & pathway systems
- https://bitbucket.org/biobakery/humann2

Faster functional profiling by avoiding translated search



Faster functional profiling by avoiding translated search



Faster functional profiling by avoiding translated search





42

HUMAnN2 accuracy (1M read mock stool metagenome)

Community Total

Alistipes onderdonkii Alistipes putredinis Alistipes shahii Bacteroides caccae Bacteroides cellulosilyticus Bacteroides dorei Bacteroides massiliensis Bacteroides ovatus Bacteroides stercoris Bacteroides thetaiotaomicron Bacteroides uniformis Bacteroides vulgatus Barnesiella intestinihominis Dialister invisus Eubacterium rectale Faecalibacterium prausnitzii Parabacteroides distasonis Parabacteroides merdae Prevotella copri Ruminococcus bromii



- 20 common gut bugs (even)
- 1M 100-nt reads
- Computed expected UniRef50 abundances from genome annotations
- Ran reads through HUMAnN2
- Compared expected and observed profiles
- Strong agreement, even for closely related species (e.g. *Bacteroides*)



	MetaPhlAn2 Prescreen	Pangenome Search	Translated Search	Total
Normal Flow	0.5 cpu-hours	0.7 cpu-hours (86% reads)	 1.7 cpu-hours (14% reads) 	2.9 cpu-hours
Translated Search Only	NA	NA	12.1 cpu-hours	12.1 cpu-hours

HUMANN2: Example combining DNA- & RNAseq from 8 healthy gut microbiomes





- Dot = functional contribution of one species
- Ribosomal & peptidoglycan transcription correlate
- Ribosome biosyn. generally "over-transcribed"
- Peptidoglycan biosyn. generally "under-transcribed"
- Not a paradox, it's consistent with the biology

https://huttenhower.sph.harvard.edu/humann2

HUMAnN2: Glycolytic processes performed by different species in Finns and Russians



https://huttenhower.sph.harvard.edu/humann2



ShortBRED

For targeted meta'omic functional profiling

The problem with short reads and regions of local homology among proteins

- Protein of interest
- Belongs to a family
- Local homology to unrelated families
- Short reads from unrelated families may map to protein of interest (spurious hits)













Jim Kaminski



Metagenome reads ShortBRED markers Translated search for high ID hits Normalize relative abundances

https://huttenhower.sph.harvard.edu/shortbred

ShortBRED Synthetic Evaluation (ABR genes)



Relative to mapping reads against full-length centroids, we are:> Substantially more accurate (fewer false positives)> Faster (reduced search space)

https://huttenhower.sph.harvard.edu/shortbred

DOC

ShortBRED: ABR in human gut metagenomes



https://huttenhower.sph.harvard.edu/shortbred

DOC



PICRUSt For predictive functional profiling

 PICRUSt: Inferring community metagenomic potential from marker gene sequencing

PICRUSt Accuracy (Spearman,r)

With Rob Knight, Rob Beiko

One can recover <u>general</u> community function with reasonable accuracy from 16S profiles.

http://picrust.github.com





Average 16S distance to nearest reference genome (NSTI)



Who is there?

What are they doing?

Sample #	1	2	3	4	5	6
Clade1	0.40	0.87	0.43	0.68	0.47	0.32
Clade1 Bug1	0.40	0.56	0.07	0.31	0.42	0.27
Clade1 Bug2	0.00	0.30	0.36	0.37	0.04	0.05
Clade2	0.60	0.13	0.57	0.32	0.53	0.68
Clade2 Bug3	0.11	0.00	0.10	0.32	0.15	0.23
Clade2 Bug4	0.49	0.13	0.47	0.00	0.39	0.45



Who is there? What are they doing? What does it all mean?

Sample #	1	2	3	4	5	6	
Profession	Student	Postdoc	Postdoc	Professor	Student	Student	
Gender	Male	Female	Female	Male	Male	Female	
Site	Oral	Gut	Oral	Gut Oral		Gut	
Clade1	0.40	0.87	0.43	0.68	0.47	0.32	
Clade1 Bug1	0.40	0.56	0.07	0.31	0.42	0.27	
Clade1 Bug2	0.00	0.30	0.36	0.37	0.04	0.05	
Clade2	0.60	0.13	0.57	0.32	0.53	0.68	
Clade2 Bug3	0.11	0.00	0.10	0.32	0.15	0.23	
Clade2 Bug4	0.49	0.13	0.47	0.00	0.39	0.45	

Properties of microbiome data

- Compositional nature (Σ = 1)
 - Abundance is relative, not absolute
- High dynamic range
- Often sparse (sample dominated by a few species)
- Noisy

DOC

• Hierarchical organization

Site	Oral	Gut	Oral	Gut	Oral	Gut
Clade1	0.40	0.87	0.43	0.68	0.47	0.32
Clade1 Bug1	0.40	0.56	0.07	0.31	0.42	0.27
Clade1 Bug2	0.00	0.30	0.36	0.37	0.04	0.05
Clade2	0.60	0.13	0.57	0.32	0.53	0.68
Clade2 Bug3	0.11	0.00	0.10	0.32	0.15	0.23
Clade2 Bug4	0.49	0.13	0.47	0.00	0.39	0.45

Properties of microbiome data

- General problem: correlate microbiome features with metadata (potentially controlling for other features)
- Intuitively summarize the results

DOC

Sample #	1	2	3	4	5	6
Profession	Student	Postdoc	Postdoc	Professor	Student	Student
Gender	Male	Female	Female	Male	Male	Female
Site	Oral	Gut	Oral	Gut	Oral	Gut
Clade1	0.40	0.87	0.43	0.68	0.47	0.32
Clade1 Bug1	0.40	0.56	0.07	0.31	0.42	0.27
Clade1 Bug2	0.00	0.30	0.36	0.37	0.04	0.05
Clade2	0.60	0.13	0.57	0.32	0.53	0.68
Clade2 Bug3	0.11	0.00	0.10	0.32	0.15	0.23
Clade2 Bug4	0.49	0.13	0.47	0.00	0.39	0.45





A more general solution for finding significant metagenomic associations in metadata-rich studies

Tim Tickle



Microbiome downstream analyses: interaction network reconstruction



It's a jungle in there – microbial interactions follow patterns from classical macro-ecology.

Mutualism

Predation

Competition



Given microbial relative abundance measurements over many samples, can we detect co-occurrence and co-exclusion relationships?

Relative abundance data poses a problem for correlating metagenomic features

	Sample1	Sample2	Sample3	sample4	sample5
Bug1	100	100	100	100	100
Bug2	1	10	100	1000	10000
	10000				/
	1000 —				
	100				
	10 -				
	1 -	sample?	ample3 sar	nplea	nes.

	Sample1	Sample2	Sample3	sample4	sample5
Bug1	0.99	0.91	0.50	0.09	0.01
Bug2	0.01	0.09	0.50	0.91	0.99



Absolute (cell) counts No bug1-bug2 correlation Relative abundance Spurious bug1-bug2 correlation (sequencing yields rel. ab.)

CCREPE: <u>Compositionality</u> <u>Corrected</u> by <u>RE</u>normalization and <u>PE</u>rmutation

Estimating a confidence interval



CCREPE: <u>Compositionality</u> <u>Corrected</u> by <u>RE</u>normalization and <u>PE</u>rmutation



- Synthetic evaluation
- Random sample feature/tables
- No built-in correlation structure

63

CCREPE: <u>Compositionality</u> <u>Corrected</u> by <u>RE</u>normalization and <u>PE</u>rmutation



"Microbial co-occurrence relationships in the human microbiome." Faust, et al. *PLoS Comp Biol*, 8:e1002606 (2012).



Who is there? What are they doing? What does it all mean?



Resources

Using tools through Galaxy

DOC

🗧 Galaxy / Huttenhov	wer Lab Analyze Data Workflow Shared Data - Help - User -	Using 0 bytes
Tools 🌣		History 🌣
Tools search tools HUTTENHOWER LAB MODULES LEFSe MetaPhIAn GraPhIAn microPITA MaAsLin PICRUSt LOAD DATA MODULE Get Data DEFAULT GALAXY MODULES Convert Formats FASTA manipulation General Galaxy tools	Thanks for visiting our lab's tools and applications page, implemented within the Galaxy web application and workflow framework. Here, we provide a number of resources for metagenomic and functional genomic analyses, intended for research and academic use. Please see the menus and folders to the left for an overview of available tools including documentation, sample data, and publications. Our lab's research interests include metagenomics and the <u>human microbiome</u> , the relationships between microbial communities and human health, microbiome systems biology. and large-scale computational methods for studying all of these areas. In addition to the tools provided here, feel free to take a look at our additional research and publications, including the <u>Sleipnir library</u> for computational functional genomics. The tools are available here without account creation. However, you are strongly invited to create an account for having access to the history, saved analyses, datasets and workflows. You can create an account and/or log in using the User menu in the top-right corner. If you have any comments, questions, or suggestions, please contact <u>Dr. Huttenhower</u> .	History O bytes Your history is empty. Click 'Get Data' on the left pane to start
<		

http://huttenhower.sph.harvard.edu/galaxy

Tutorials available online

DOC



http://huttenhower.sph.harvard.edu/biobakery (click on your tool-of-interest)

All tools are open source

DOC

≡	🖲 Bitbucket	Teams - Rej	positories 👻	Create			Q owner/repo	sitory	? ◄	@ -
	biobak	ery ry ⊠ Share				Ł Clone →	🅽 Branch 🚹	Pull request	•••	•
Ove	erview Source	Commits	Branches	Pull requests	Issues 1	Wiki Do	ownloads			\$
Ho	ome						Clone wiki 👻	Edit Crea	ate H	istory
	Huttenhower Lab Tools Welcome to the official Huttenhower Tutorials wiki. The wiki follows through the computational tools currently being used by our lab, which is also publicly available. The tools can be divided under three categories: Composition Analysis These tools can determine the composition in terms of (i) microbial species and their associated abundances (MetaPhlAn) or (ii) genes and associated genes (HUMAnN) in the dataset. Please click on the links below for detailed tutorials: PublePhlan • PhyloPhlAn • Nicrobial species and associated genes and abundances of proteins of phylogenetic trees • PhyloPhlan • Nicrobial species and abundances • Reconstruction of phylogenetic trees • PhyloPhlan • Reconstruction of phylogenetic trees • Bundance of proteins of interest of genetic data • Dundance of proteins of interest of genetic data									

http://bitbucket.org/biobakery/biobakery

The bioBakery Virtual Machine

DOC

https://bitbucket.org/biobakery/biobakery/wiki/biobakery_wiki



Ubuntu base image preloaded and configured to run all Huttenhower lab tools; one click up-and-running via Vagrant

Thank you! DOC

Xochitl

Morgan

Gholamali

Rahnavard

Boyu

Ren



Curtis Huttenhower



Aleksandar Kostic



Tiffany Hsu



Casey DuLong



Luc Bijnens

Tommi Vatanen



Emma

Schwager

Minah

lqbal

Melanie

Schirmer

Jim Kaminski







Levi Waldron Nicola Segata



Siyuan

Ma



Wendy Garrett **Michelle Rooks**



Kat Huang



Ramnik Xavier MGH 1811 Harry Sokol Dan Knights Moran Yassour



W **Rob Beiko** Morgan Langille











Mark Silverberg Rob Knight Boyko Kabakchiev Greg Caporaso Andrea Tyler Jesse Zaneveld

Bruce Sands







Sahar Abubucker

Mathangi Thiagarajan

Beltran Rodriguez-Mueller

Makedonka Mitreva

Yuzhen Ye

Mihai Pop

Larry Forney

Brandi Cantarel Alyx Schubert

George Weingart



 \otimes Ayshwarya









Bruce Birren Mark Daly Dovle Ward Ashlee Earl





Ruth Lev

Omry Koren

Owen White

Lita Proctor

Joe Petrosino

George Weinstock Karen Nelson

Erica Sodergren

Anthony Fodor

Marty Blaser

Jacques Ravel





Sirota-Madi Subramanian

Lauren

Mclver

Alexandra

Koji

Yasuda

Andy Shi