# ONE CODEX

# A platform for highly accurate, reproducible metagenomics

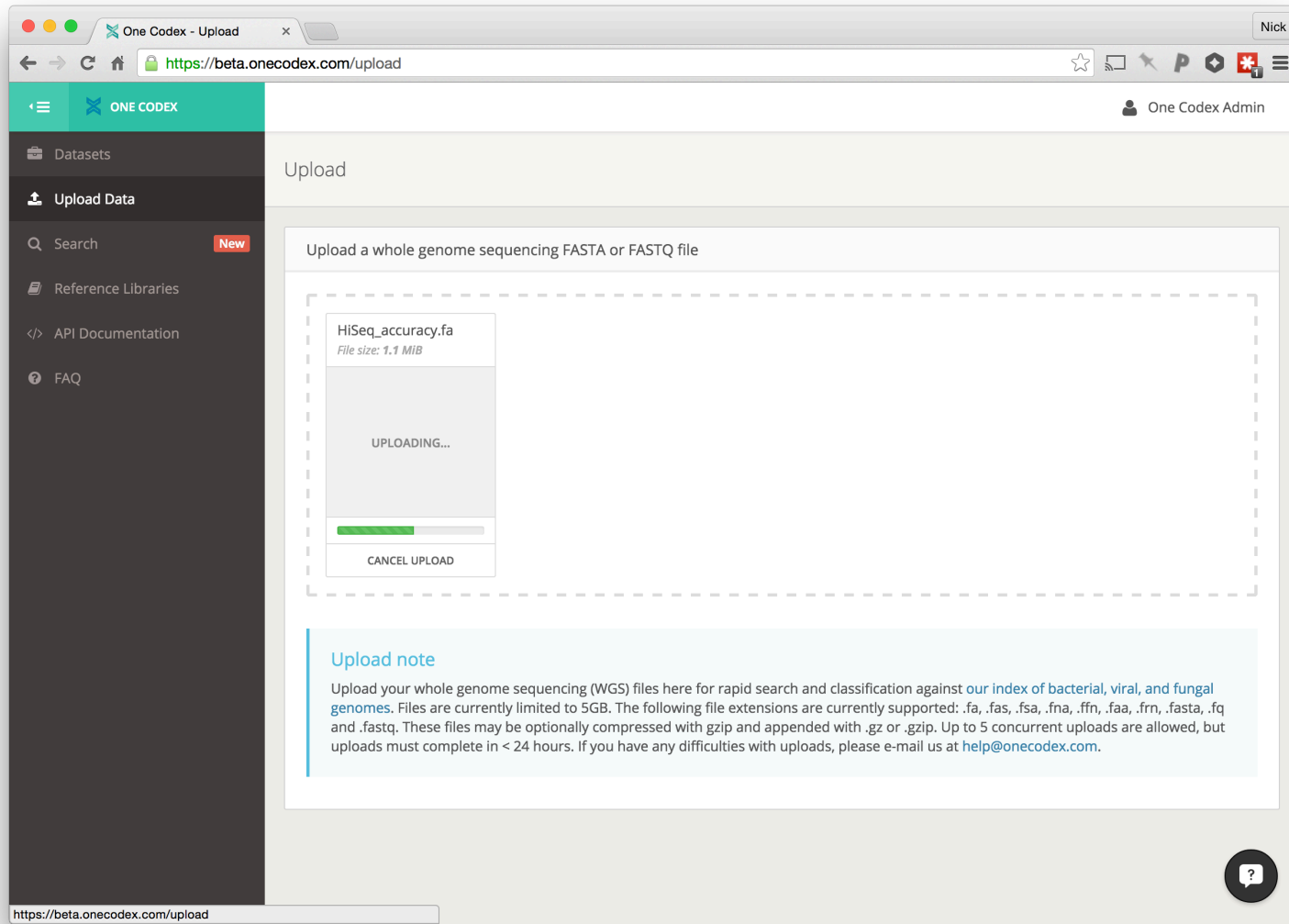*MetaSUB Meeting – New York Genome Center*

*June 20, 2015*

# Quick overview

- One Codex
  - Microbial genomics software and data platform
  - Based in San Francisco

- How can we make better applied genomics software?
  - Usable + intuitive
  - Open + extensible (e.g., APIs)
  - Move beyond analyzing a single sample (enable "databanking")

# Towards applied microbial genomics

*Technology should "just work"…*
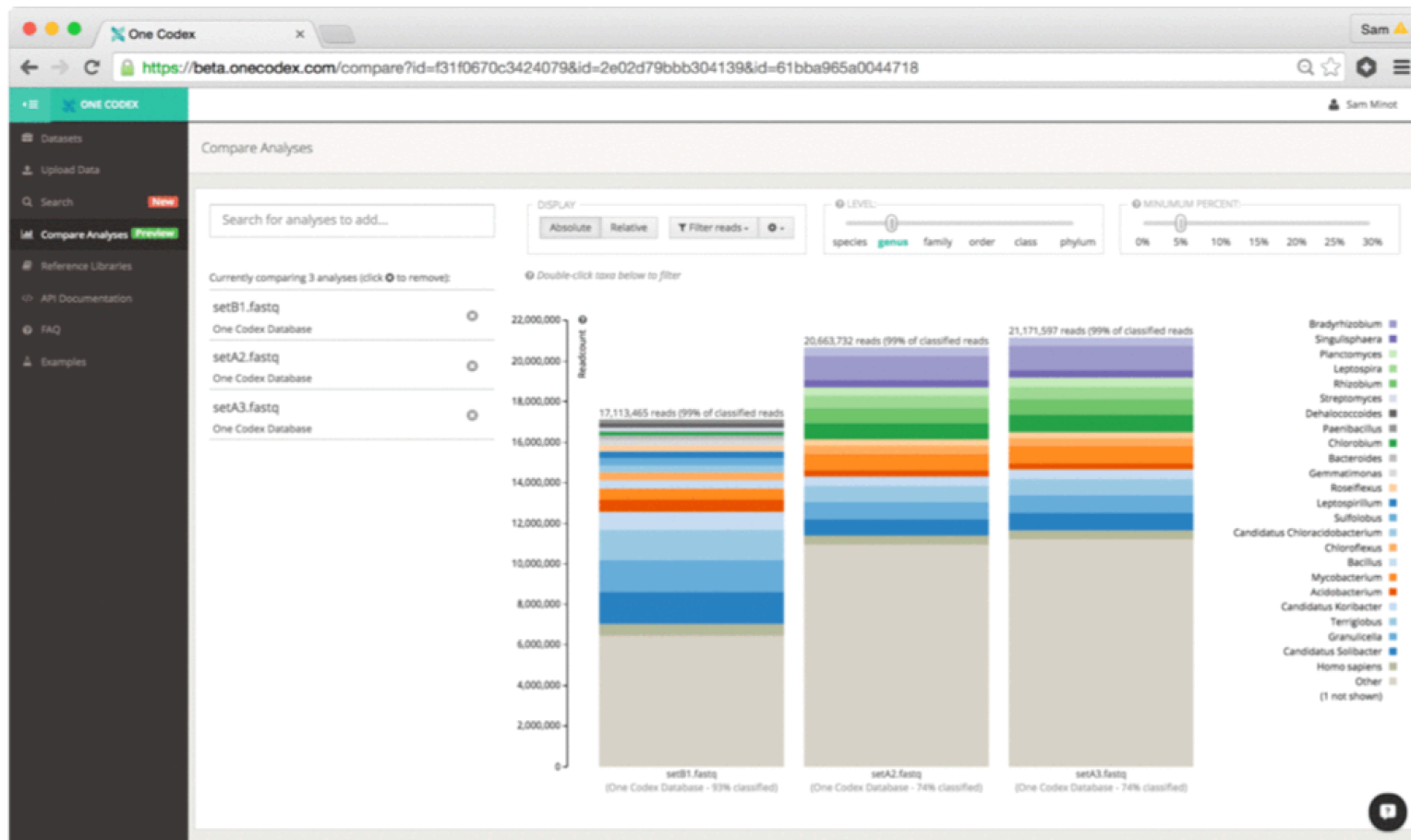
ONE CODEX

# Towards applied microbial genomics

*… if we're going to sequence everything, we should be able to easily index, explore, and compare samples*

# Towards applied microbial genomics

*… if we're going to sequence everything, we should be able to easily index, explore, and compare samples*

# Towards applied microbial genomics

*… if we're going to sequence everything, we should be able to easily index, explore, and compare samples*

ONE CODEX

# Algorithms/methods at One Codex

- "Exact alignment" or "alignment-free" *k*-mer based methods

- Scale well to huge reference sets (have 35k today)

- Dive into two today:

  a. Metagenomic classification

  b. Strain-typing

ONE CODEX

# #1 Build a reference library



*Reference Genomes*

- Our latest database has ~35K bacterial, viral, fungal, and archaeal genomes

- Bigger tends to be better

ONE CODEX

# #2 Enumerate all possible *k*-mers

Reference Genomes

k-mers

ONE CODEX

# #3 Associate *k*-mers and tax IDs

Reference Genomes                    *k*-mers

ONE CODEX

# #3 Associate *k*-mers and tax IDs

Select all the unique *k*-mers

Reference Genomes

*k*-mers

ONE CODEX

# #3 Associate *k*-mers and tax IDs

Select all the unique *k*-mers

**Assign parent tax IDs**

E. Coli (562)

Salmonella (590)

Entero-bactericeae (543)

Reference Genomes

*k*-mers

ONE CODEX

# #3 Associate *k*-mers and tax IDs

Select all the unique *k*-mers

Assign parent tax IDs

**Assign LCA for shared *k*-mers**

E. Coli (562)

Salmonella

K. Mil

Enterobacteriaceae

*Reference Genomes*

*k-mers*

ONE CODEX

# #4 Lookup all *k*-mers for an input

```
@Sample.FASTQ length=151
TAGAGAATGTTAAGTCAGTTCCAGGACCAGTTAGGCGCACTTTATCTGTTTTATT...
+
BBBBBFFFFFFFFGGGFGGBGGGGHGHHCGHHHGGHAGFHHDFHFGHHHFHHFGGE...
```

# #4 Lookup all *k*-mers for an input

```
@Sample.FASTQ length=151
TAGAGAATGTTAAGTCAGTTCCAGGACCAGTTAGGCGCACTTTATCTGTTTTATT...
```

Firmicutes (1239)

Bacillus cereus group (86661)

Bacillus cereus (1396)

# #5 Classify the input read

`@Sample.FASTQ length=151`
`TAGAGAATGTTAAGTCAGTTCCAGGACCAGTTAGGCGCACTTTATCTGTTTTATT...`

Firmicutes (1239)

Bacillus cereus group (86661)

Bacillus cereus (1396)

*Classify result:*
*Highest weighted*
*root-to-leaf path*

ONE CODEX

# Metagenomic classification summary

**BUILD**

1.  Build a library of reference genomes

2.  Enumerate all of the $k$-mers

3.  Build a database of $k$-mers and associated taxonomy IDs

**CLASSIFY**

4.  Look up all possible $k$-mers in a sample

5.  Classify each read independently based on these $k$-mer "hits"

ONE CODEX

# Metagenomic classification: accuracy

*The One Codex platform delivers a comprehensive database – enabling substantial accuracy improvements*

| Classifier | HiSeq – Precision | HiSeq – Sensitivity | MiSeq – Precision | MiSeq – Sensitivity |
|---|---|---|---|---|
| PhymmBL | 79.1 | 79.1 | 76.2 | 76.2 |
| PhymmBL (conf. > 0.65) | 99.1 | 73.9 | 92.5 | 73.0 |
| Megablast | 99.1 | 79.0 | 92.4 | 75.7 |
| Naïve Bayes | 82.3 | 82.3 | 77.8 | 77.8 |
| Kraken | 99.2 | 77.1 | 94.7 | 73.5 |
| **One Codex** | **99.5** | **96.4** | **97.8** | **90.4** |

Comparison of genus-level performance values adapted from Wood and Salzburg, 2014

ONE CODEX

# Metagenomic classification: results

# Metagenomic classification: results

# B) Strain-typing

- Won CDC "No Petri Dish" Challenge

- Works for isolates and mixed samples

- 99+% "serotype-level" accuracy

- Higher resolution than MLST at 10-20x lower coverage

ONE CODEX

# B) Strain-typing

ONE CODEX

# B) Strain-typing: accuracy

**Strain-typing accuracy for 722 E. coli isolate datasets**

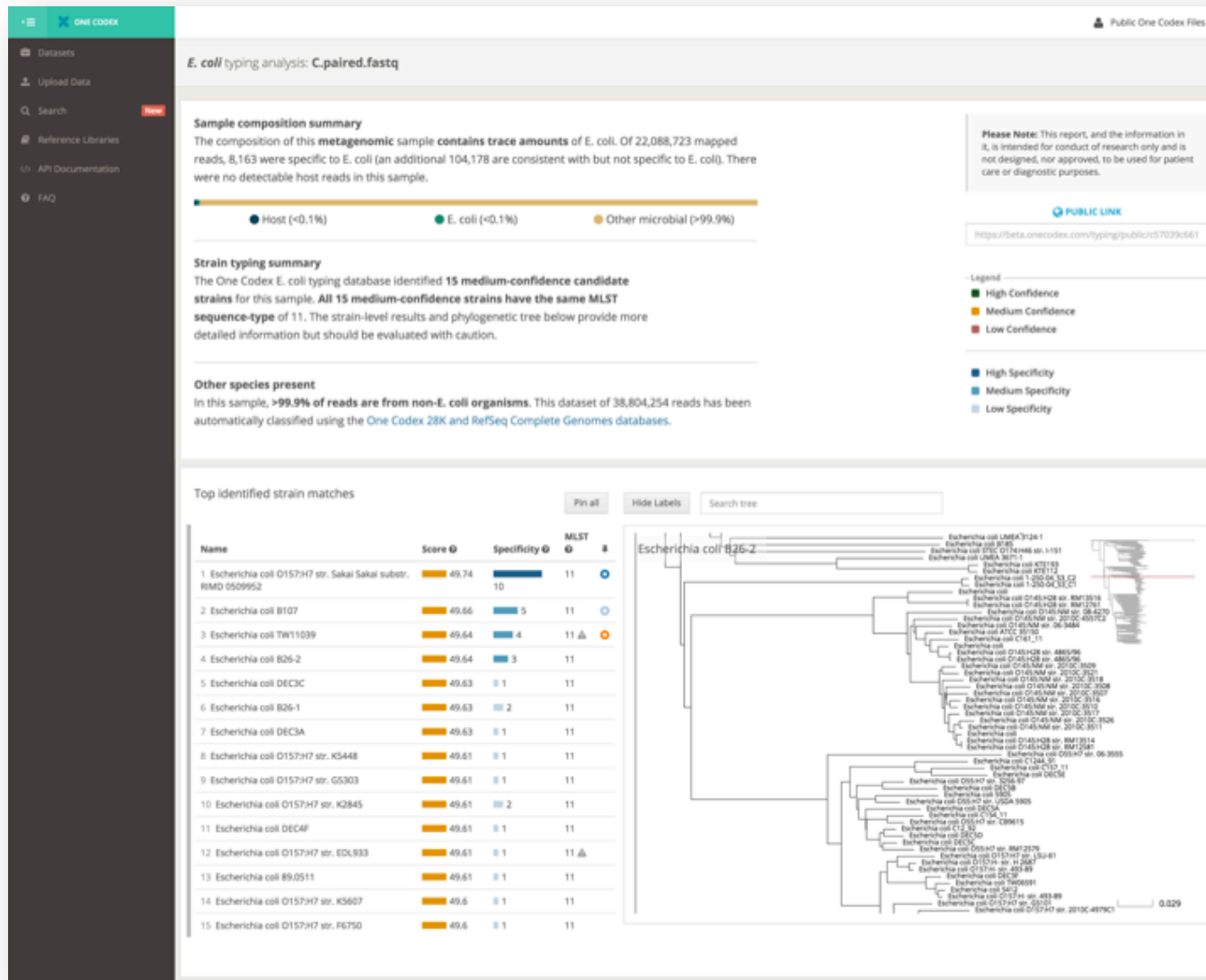| Accuracy level | Number of correctly identified strains | Percentage |
|---|---|---|
| **Serotype-level accuracy** [1] | 718/722 | 99.5% |
| **Outbreak-level accuracy** [2] | 713/722 | 98.7% |

**STEC identification of 3 STEC strains in 3 stool samples**

| Spike % of total | Spiked reads | *E. coli* genome coverage | **Outbreak-level** [2] accuracy | **Serotype-level** [1] accuracy |
|---|---|---|---|---|
| **10%** | 1M–2M | 40–80x | **100%** | **100%** |
| **5%** | 0.5–1M | 20–40x | **100%** | **100%** |
| **1%** | 100K–200K | 4–6x | **100%** | **100%** |
| **0.1%** | 10K–20K | 0.4–1.2x | **100%** | **100%** |
| **0.05%** | 5K–10K | 0.2–0.5x | **89%** (8/9 samples) | **100%** |
| Control | None | | No STEC detected | |
| | | **Overall** | **98% (44/45)** | **100% (45/45)** |

*Tested strains: O157:H7 str. F8092B, O157:H7 str. Sakai, and O104:H4 str. TY-2482. Stool samples from Human Microbiome Project.*

[1] Identified strain within 0.06 genetic distance (median within-serotype distance) from known spike-in
[2] Identified strain within 95% confidence interval for intra-serotype genetic distance

ONE CODEX

# Questions?

Open beta: https://beta.onecodex.com

Contact: nick@onecodex.com