

Microbial Typing by Machine Learned DNA Melt Signatures

**Nadya Andini¹, Bo Wang², Pornpat Athamanolap³, Justin Hardick⁴, Billie J. Masek⁵, Simone Thair¹,
Annie Hu¹, Gideon Avornu⁵, Stephen Peterson⁵, Steven Cogill¹, Richard E. Rothman^{4,5}, Karen C.
Carroll⁶, Charlotte A. Gaydos^{4,5}, Tza-Huei Wang^{3,7}, Serafim Batzoglou², Samuel Yang^{1,*}**

¹Emergency Medicine, Stanford University, Stanford, California, 94305, USA.

²Computer Science, Stanford University, Stanford, California, 94305, USA.

³Biomedical Engineering, The Johns Hopkins University, Baltimore, Maryland, 21218, USA.

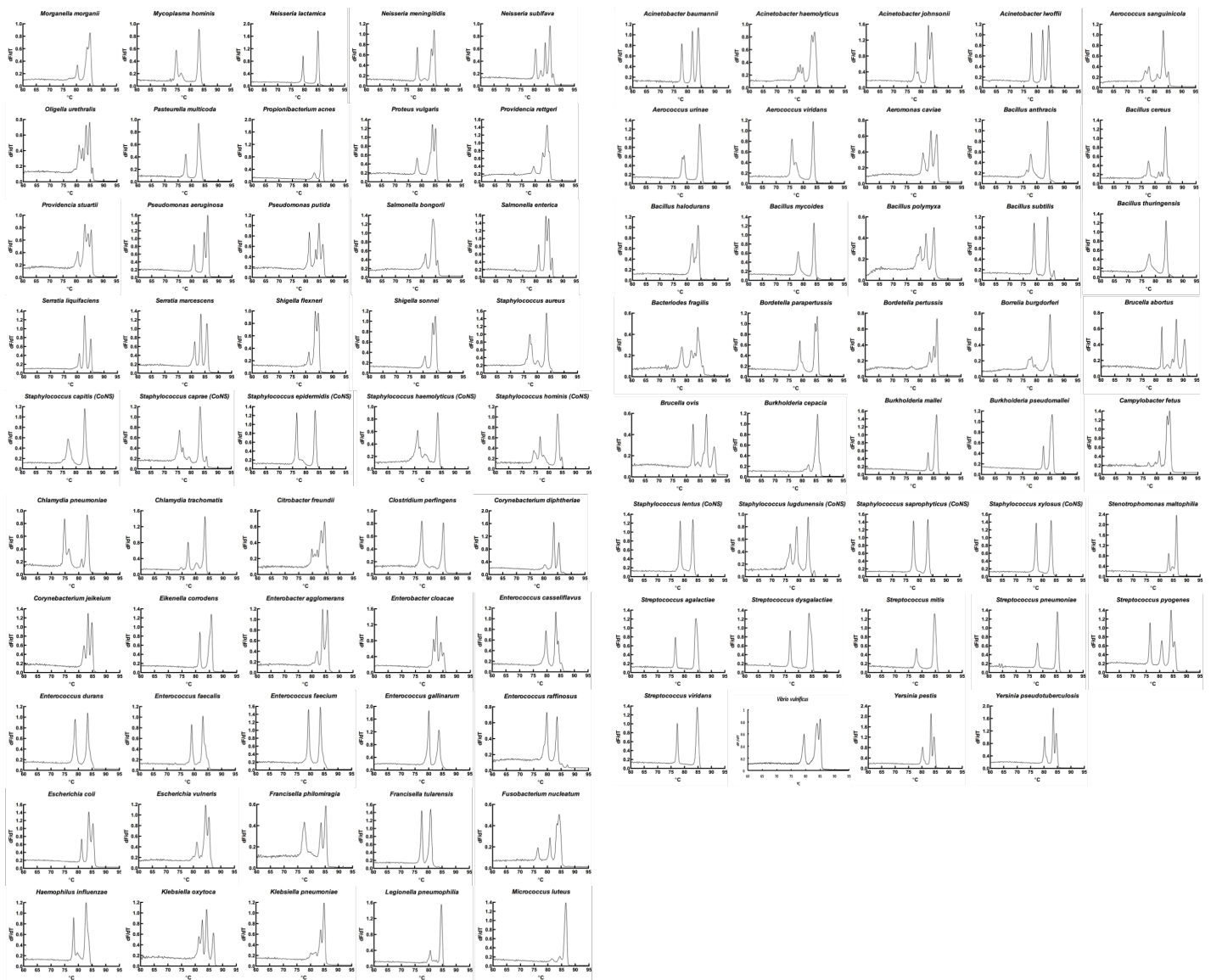
⁴Infectious Disease, Medicine, The Johns Hopkins University, Baltimore, Maryland, 21218, USA.

⁵Emergency Medicine, The Johns Hopkins University, Baltimore, Maryland, 21218, USA.

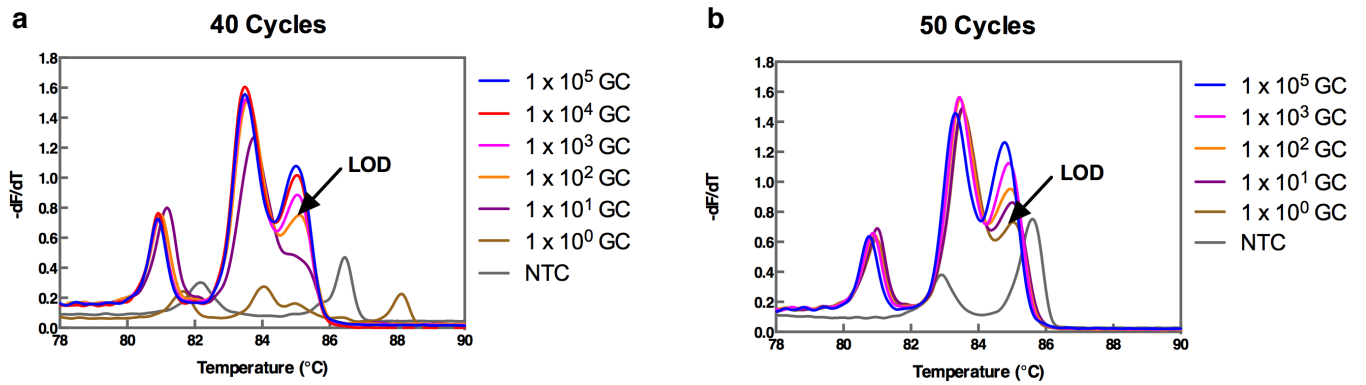
⁶Medical Microbiology, Pathology, The Johns Hopkins University, Baltimore, Maryland, 21218, USA.

⁷Mechanical Engineering, The Johns Hopkins University, Baltimore, Maryland, 21218, USA.

*corresponding.syang5@stanford.edu



Supplementary Fig. S1. Individual derivative melt curves of 89 reference bacterial species.



Supplementary Fig. S2. The limit of detection analysis of ITS PCR HRM. Serially diluted *E. coli* genomic DNA calculated based on its genome copies (GC) was amplified in a 40 (a) and a 50 (b)-cycle PCR targeting the ITS region. The PCR was immediately followed by HRM to produce corresponding derivative melt curves. The limit of detection (LOD) was determined to be the concentration where melt curve profile was maintained (arrows).

Supplementary Table S1. List of 89 reference bacterial species in the database.

Bacterial Species in Database			
<i>Acinetobacter baumannii</i>	<i>Burkholderia pseudomallei</i>	<i>Haemophilus influenzae</i>	<i>Shigella flexneri</i>
<i>Acinetobacter haemolyticus</i>	<i>Campylobacter fetus</i>	<i>Klebsiella oxytoca</i>	<i>Shigella sonnei</i>
<i>Acinetobacter johnsonii</i>	<i>Chlamydia pneumoniae</i>	<i>Klebsiella pneumoniae</i>	<i>Staphylococcus aureus</i>
<i>Acinetobacter lwoffii</i>	<i>Chlamydia trachomatis</i>	<i>Legionella pneumophila</i>	<i>Staphylococcus capitis (CoNS)</i>
<i>Aerococcus sanguinicola</i>	<i>Citrobacter freundii</i>	<i>Micrococcus luteus</i>	<i>Staphylococcus caprae (CoNS)</i>
<i>Aerococcus urinae</i>	<i>Clostridium perfringens</i>	<i>Morganella morganii</i>	<i>Staphylococcus epidermidis (CoNS)</i>
<i>Aerococcus viridans</i>	<i>Corynebacterium diphtheriae</i>	<i>Mycoplasma hominis</i>	<i>Staphylococcus haemolyticus (CoNS)</i>
<i>Aeromonas caviae</i>	<i>Corynebacterium jeikeium</i>	<i>Neisseria lactamica</i>	<i>Staphylococcus hominis (CoNS)</i>
<i>Bacillus anthracis (2 strains)</i>	<i>Eikenella corrodens</i>	<i>Neisseria meningitidis</i>	<i>Staphylococcus lentus (CoNS)</i>
<i>Bacillus cereus</i>	<i>Enterobacter agglomerans</i>	<i>Neisseria subflava</i>	<i>Staphylococcus lugdunensis (CoNS)</i>
<i>Bacillus halodurans</i>	<i>Enterobacter cloacae</i>	<i>Oligella urethralis</i>	<i>Staphylococcus saprophyticus (CoNS)</i>
<i>Bacillus mycoides</i>	<i>Enterococcus casseliflavus</i>	<i>Pasteurella multocida</i>	<i>Staphylococcus xylosus (CoNS)</i>
<i>Bacillus polymyxa</i>	<i>Enterococcus durans</i>	<i>Propionibacterium acnes</i>	<i>Stenotrophomonas maltophilia</i>
<i>Bacillus subtilis</i>	<i>Enterococcus faecalis</i>	<i>Proteus vulgaris</i>	<i>Streptococcus agalactiae</i>
<i>Bacillus thuringiensis</i>	<i>Enterococcus faecium</i>	<i>Providencia rettgeri</i>	<i>Streptococcus dysgalactiae</i>
<i>Bacteriodes fragilis</i>	<i>Enterococcus gallinarum</i>	<i>Providencia stuartii</i>	<i>Streptococcus parasanguinis</i>
<i>Bordetella parapertussis</i>	<i>Enterococcus raffinosus</i>	<i>Pseudomonas aeruginosa</i>	<i>Streptococcus pneumoniae (5 strains)</i>
<i>Bordetella pertussis</i>	<i>Escherichia coli</i>	<i>Pseudomonas putida</i>	<i>Streptococcus pyogenes</i>
<i>Borrelia burgdorferi</i>	<i>Escherichia vulneris</i>	<i>Salmonella bongorii</i>	<i>Streptococcus anginosus</i>
<i>Brucella abortus</i>	<i>Francisella philomiragia</i>	<i>Salmonella enterica Enteritidis</i>	<i>Vibrio fluvialis</i>
<i>Brucella ovis</i>	<i>Francisella tularensis</i>	<i>Serratia liquifaciens</i>	<i>Yersinia pestis (2 strains)</i>
<i>Burkholderia cepacia</i>	<i>Fusobacterium nucleatum</i>	<i>Serratia marcescens</i>	<i>Yersinia pseudotuberculosis</i>
<i>Burkholderia mallei</i>			

SUPPLEMENTARY METHODS

1. Naïve Bayes

In this section, we present details about the proposed adaptive Naïve Bayes algorithm. Given C species in the reference dataset, and for the i -th species C_i , we have N_i number of training samples. For any new unknown test sample x , we aim to calculate the posteriori probability via Bayes' theorem:

$$p(C_k|x) = \frac{p(C_k)p(x|C_k)}{p(x)}$$

where $p(C_k)$ is the prior for the k -th species, and $p(x|C_k)$ is the likelihood function given all the training samples in the k -th species.

The prior information is assumed to be homogeneous:

$$p(C_k) = \frac{1}{C}$$

The likelihood function is calculated with a Gaussian distribution:

$$p(x|C_k) = \sum_{x'_j \in C_k} \frac{1}{\sqrt{2\pi\alpha}} \exp\left(-\frac{D(x, x'_j)}{2\alpha}\right).$$

The essence in the algorithm lies in the way we calculate the distance $D(x, x'_j)$. This measures the similarity between curve shapes for the test sample and training reference.

Assume for a test species, denoted as $S_t = \{S_t^1, S_t^2, \dots, S_t^m\}$ where m is the number of replicates in this species. We want to achieve a consensus prediction of whether this species falls into any species category from the reference panel. We assume each replicate of same importance, so we just average the final posteriori probability of each replicate to obtain the prediction for the test species:

$$p(C_k|S_t) = \frac{1}{m} \sum_j^m p(C_k|S_t^j)$$

2. Curve Similarity Calculation

There are three steps in the calculation of curve similarity. First, we align each curve according to the temperature of 53°C. This guarantees each curve is well-aligned and thus high accuracy in the following

curve similarity calculation. Second, we apply Hilbert Transformation on the curves. Hilbert transformation is a convolution process on the curve:

$$H(f)(t) = \frac{1}{\pi} \int_{-\infty}^{\infty} \frac{f(\tau)}{t - \tau} d\tau$$

where $f(t)$ denotes the curve we have. The output of Hilbert transformation is a complex function where the real part is the original input and the complex part denotes the transformed domain. We calculate the distance between two curves by combining the two parts as follows:

$$D(f, g) = \sum_t \|\mathbf{real}(H(f)(t)) - \mathbf{real}(H(g)(t))\|^2 + \sum_t \|\mathbf{complex}(H(f)(t)) - \mathbf{complex}(H(g)(t))\|^2$$

where f and g represent two curves.

3. Details in predicting out-of-reference samples

To distinguish whether the test target belongs to any species in the reference panel, we adapt the original Naïve Bayes to accommodate the prediction of out-of-reference samples. Assume for a test species, denoted as $S_t = \{S_t^1, S_t^2, \dots, S_t^m\}$ where m is the number of replicates in this species. First, for each replicate, we assign a prior probability to be out-of-reference sample by looking at the curve region between temperature 52.5°C to 53.5°C. This would give us some knowledge about whether this replicate is an outlier because most of outlier curves will generate some unusual peak curves in this temperature region. Further, when we apply Naïve Bayes, we assign the posteriori probability to be out-of-reference by adding a gated function that if the following quantile is below some threshold:

$$P(S_t \in C_0) = I\{\max_k p(S_t | C_k) < \theta\}.$$

we set $\theta = 0.3$ in our experiments.