

**LCA (Least Common
Ancestor)**

Least-common ancestor (LCA)

- Problem: we need a computationally efficient way to analyze millions of reads with respect to their taxonomic origin.
- Idea: analyze small subsequences instead of complete sequences; combine information from these small subsequences in a second step.
- Try to make the most specific assignment for each subsquence
- Definition: a k-mer is a sequence of length k
Example: ACGTT consists of the 3 3-mers ACG, CGT, GTT

Wood and Salzberg *Genome Biology* 2014, **15**:R46
<http://genomebiology.com/2014/15/3/R46>



METHOD

Open Access

Kraken: ultrafast metagenomic sequence classification using exact alignments

Derrick E Wood^{1,2*} and Steven L Salzberg^{2,3}

Wood et al. *Genome Biology* (2019) 20:257
<https://doi.org/10.1186/s13059-019-1891-0>

Genome Biology

SHORT REPORT

Open Access

Improved metagenomic analysis with Kraken 2

Derrick E. Wood^{1,2}, Jennifer Lu^{2,3} and Ben Langmead^{1,2*}



k-mers

- k-mer = Sequence of length k

0	1	2	3	4	5	6	7	8
T	G	A	T	A	C	G	A	A

T	G	A	T	A				
	G	A	T	A	C			
		A	T	A	C	G		
			T	A	C	G	A	
				A	C	G	A	A

LCA: 3 example genomes

Species / strain

E. coli K12

E. coli O157:H7

M. tuberculosis strain H37Rv

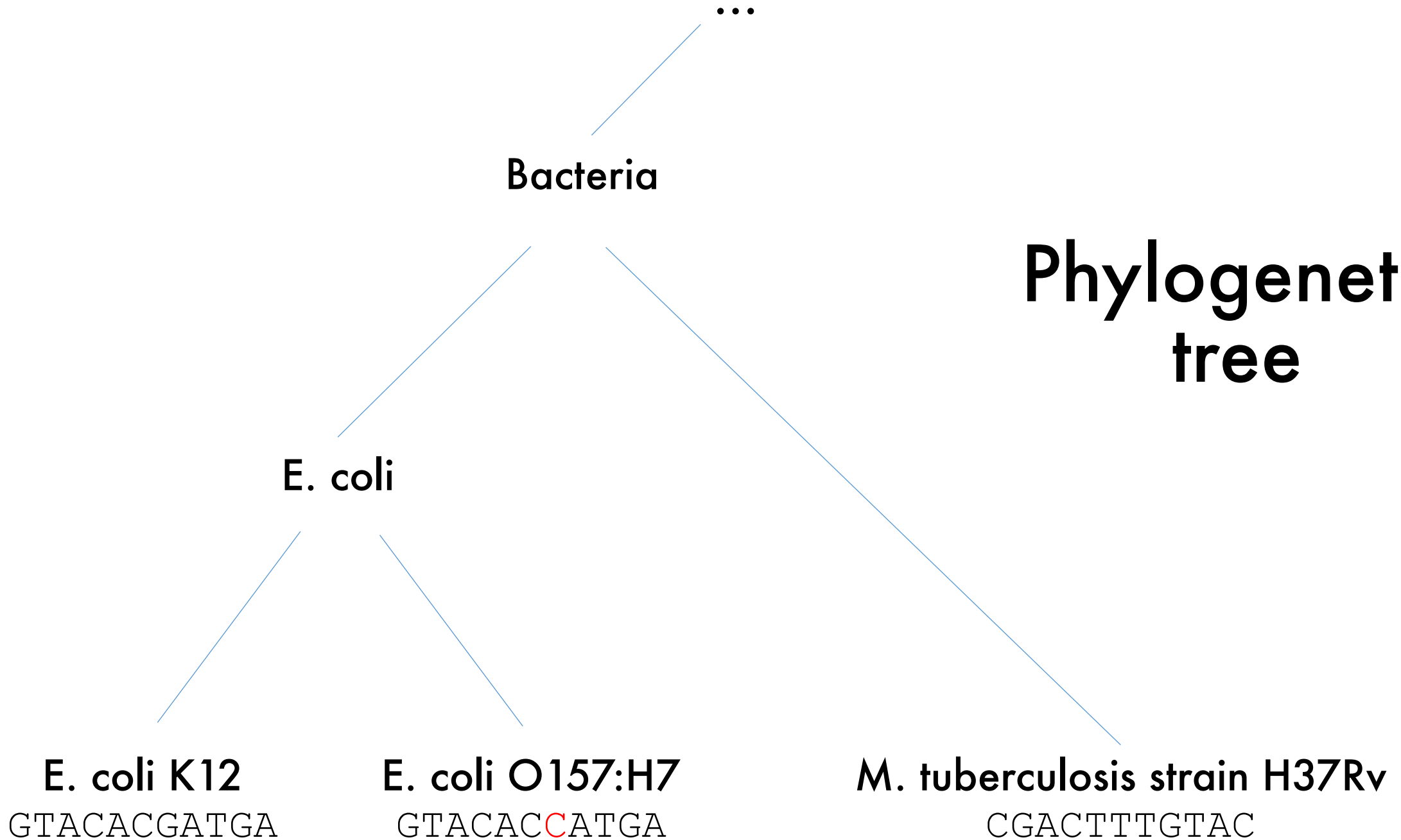
Reference genome

GTACACGATGA

GTACACCATGA

CGACTTTGTAC

Phylogenetic tree



LCA: Building the database

Split the input genomes into k-mers (e.g. 4-mers) and determine the phylogenetic placement of each k-mer.

(I.e. determine the lowest node in the phylogenetic tree that fulfils the condition that all genomes that contain the k-mer are below the chosen node).

Species / strain

E. coli K12

E. coli O157:H7

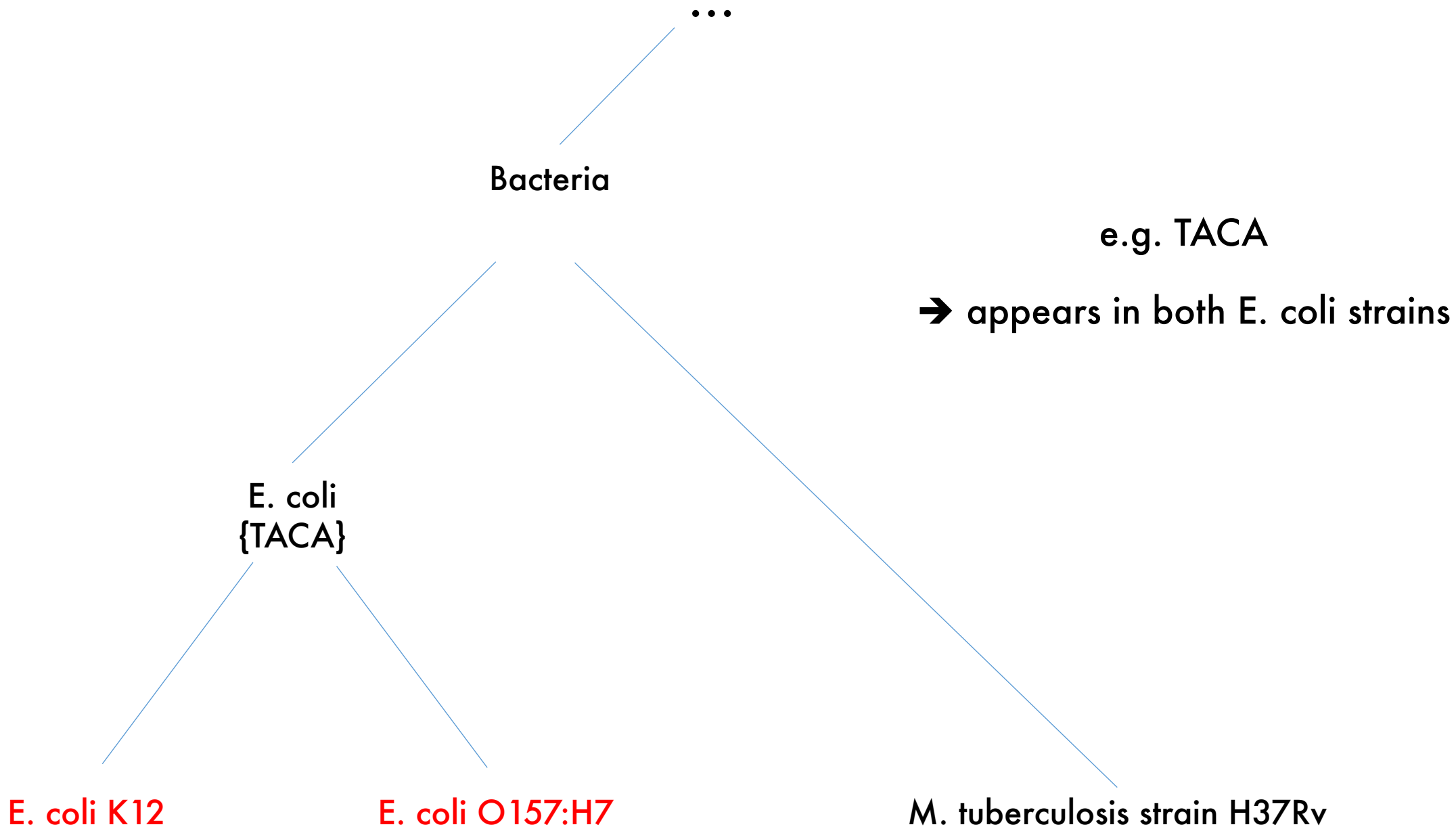
M. tuberculosis strain H37Rv

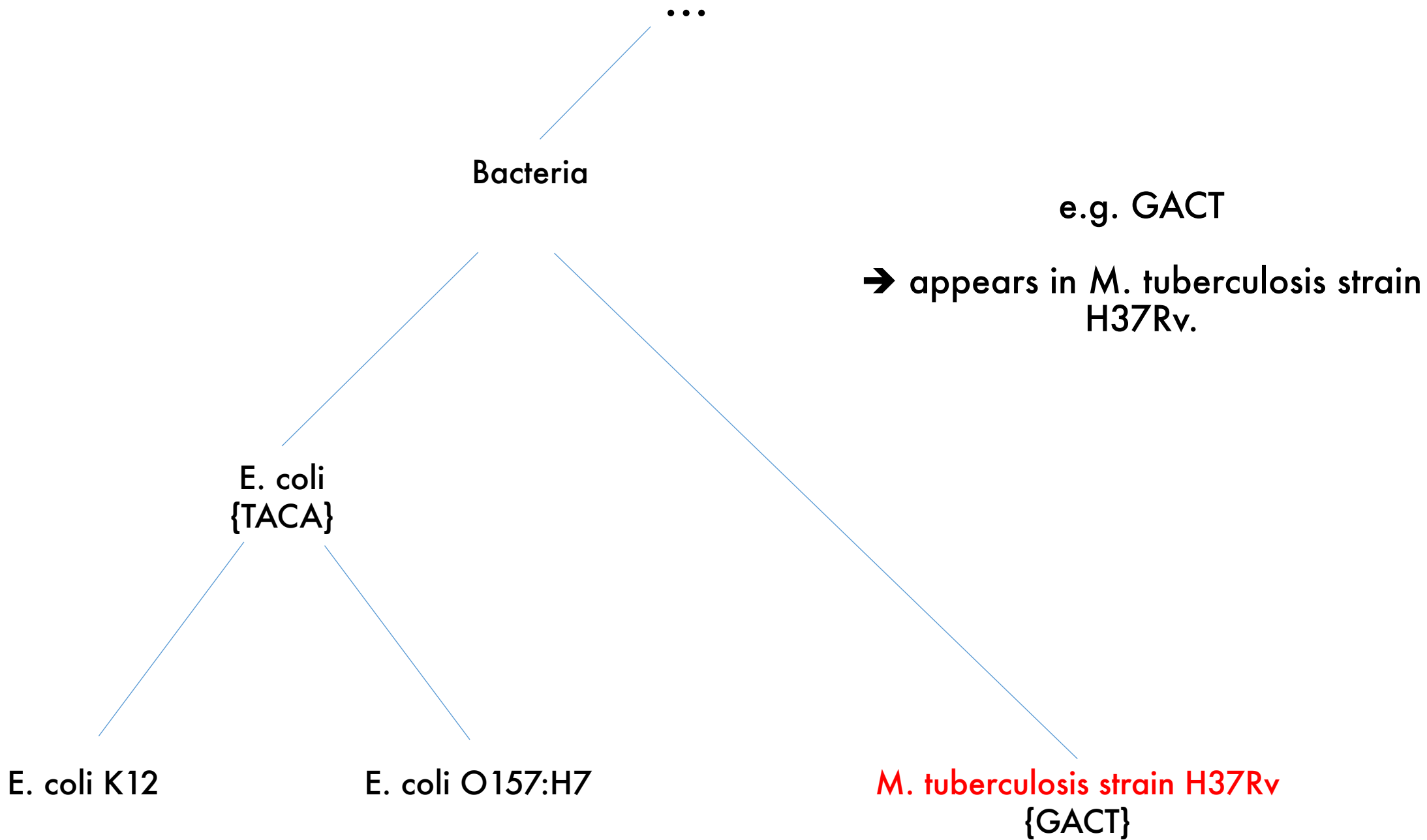
Reference genome 4-mers

{GTAC, TACA, ACAC, CACG, ACGA, CGAT, GATG, ATGA}

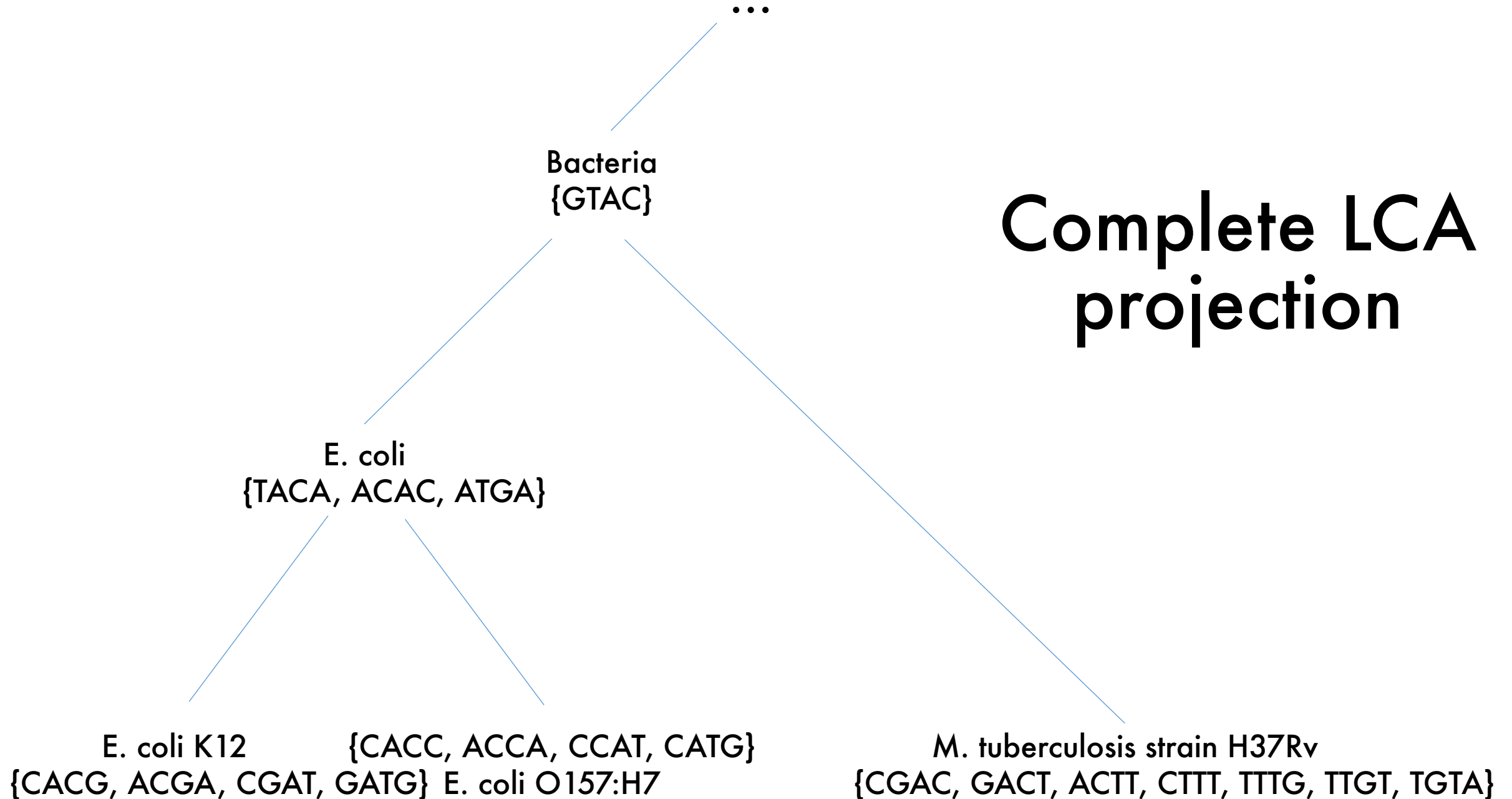
{GTAC, TACA, ACAC, CACC, ACCA, CCAT, CATG, ATGA}

{CGAC, GACT, ACTT, CTTT, TTTG, TTGT, TGTA, GTAC}





Complete LCA projection



...

Bacteria
{GTAC}

Read GACTT -> {GACT, ACTT}

E. coli
{TACA, ACAC, ATGA}

E. coli K12 {CACG, ACGA, CGAT, GATG} E. coli O157:H7 {CACC, ACCA, CCAT, CATG}

M. tuberculosis strain H37Rv {CGAC, GACT, ACTT, CTTT, TTG, TTGT, TGTA}

...

Bacteria
{GTAC}

Read GTACA -> {GTAC, TACA}

E. coli
{TACA, ACAC, ATGA}

E. coli K12 {CACG, ACGA, CGAT, GATG} E. coli O157:H7 {CACC, ACCA, CCAT, CATG}

M. tuberculosis strain H37Rv {CGAC, GACT, ACTT, CTTT, TTTG, TTGT, TGTA}

...

Bacteria
{GTAC}

Read ACGAT -> {ACGA, CGAT}

E. coli
{TACA, ACAC, ATGA}

E. coli K12

{CACC, ACCA, CCAT, CATG}

{CACG, ACGA, CGAT, GATG} E. coli O157:H7

M. tuberculosis strain H37Rv
{CGAC, GACT, ACTT, CTTT, TTTG, TTGT, TGTA}

Task [15 minutes]

Species / strain

Homo sapiens

Pan troglodytes (chimp)

E. coli O157:H7

M. tuberculosis strain H37Rv

Reference genome

CACGACGTACG

CATGACGTCCG

GTACACCATGA

CGACTTTGTAC

- What is the LCA of these species?
- Build an LCA classification tree (6-mers)
- Where do the LCA hits of the k-mers in the read `ACGACGTC` localize in the tree?
- Where would we assign the read?
- Does this sequence exist in any of our reference genomes?
- Which conclusions would we draw from that?