# Genome annotation

# Sequences

- Sequencing

- Mapping

- Assembly

# Annotation

- Identifying the locations of genes and all of the coding regions

- Determining their functions

# Annotation steps

- Identify non-coding part of genome
- Identify genome elements
- Assign functional information

# Genome elements

- Coding
  - genes
  - Transcription and translation

- Non-coding
  - Structural DNA (not transcribed)
  - Functional RNA (not translated)
  - Introns (removed before translation)

# Non-coding

- Structural
  - Telomeres, centromeres, repetetives
- Functional
  - tRNA, rRNA
- Introns
  - Removed from mRNA

# Structural annotation

- Open Reading Frames

- Gene structure

- Coding regions

- Regulatory motifs

# ORFs

- Part of genome between start and stop codons
- 6 frames for one sequence
- 1, 2, 3, -1, -2, -3

# Search potential genes

- BLAST+, HMM search, KRAKEN
- Comparing with known set of genes
- Score the similarity

# Search potential product

- Compare translated product with known database

- DIAMOND

# Identification ab initio

- Search for patterns
- AI to derive function
- GLIMMER, GeneScan

# Databases

- ENCODE
- Entrez Gene
- Ensembl
- GENCODE
- Gene Ontology Consortium
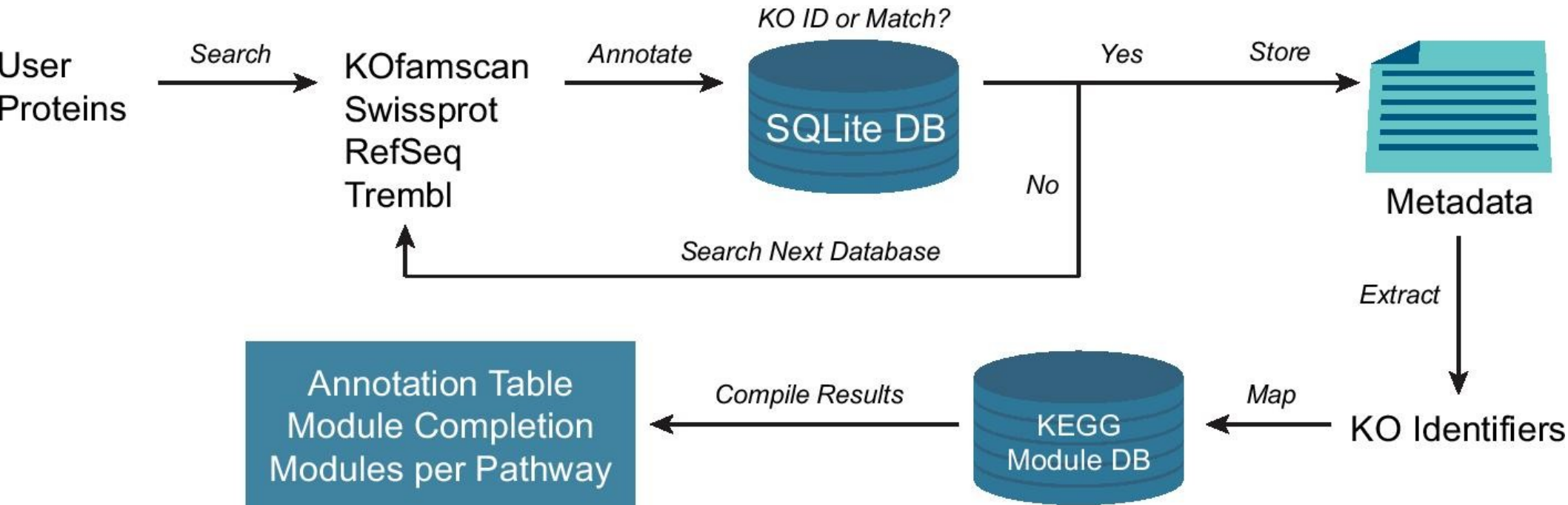- RefSeq
- Uniprot
- Vega
- ...

# Useful tools

- Prokka

- MicrobeAnnotator

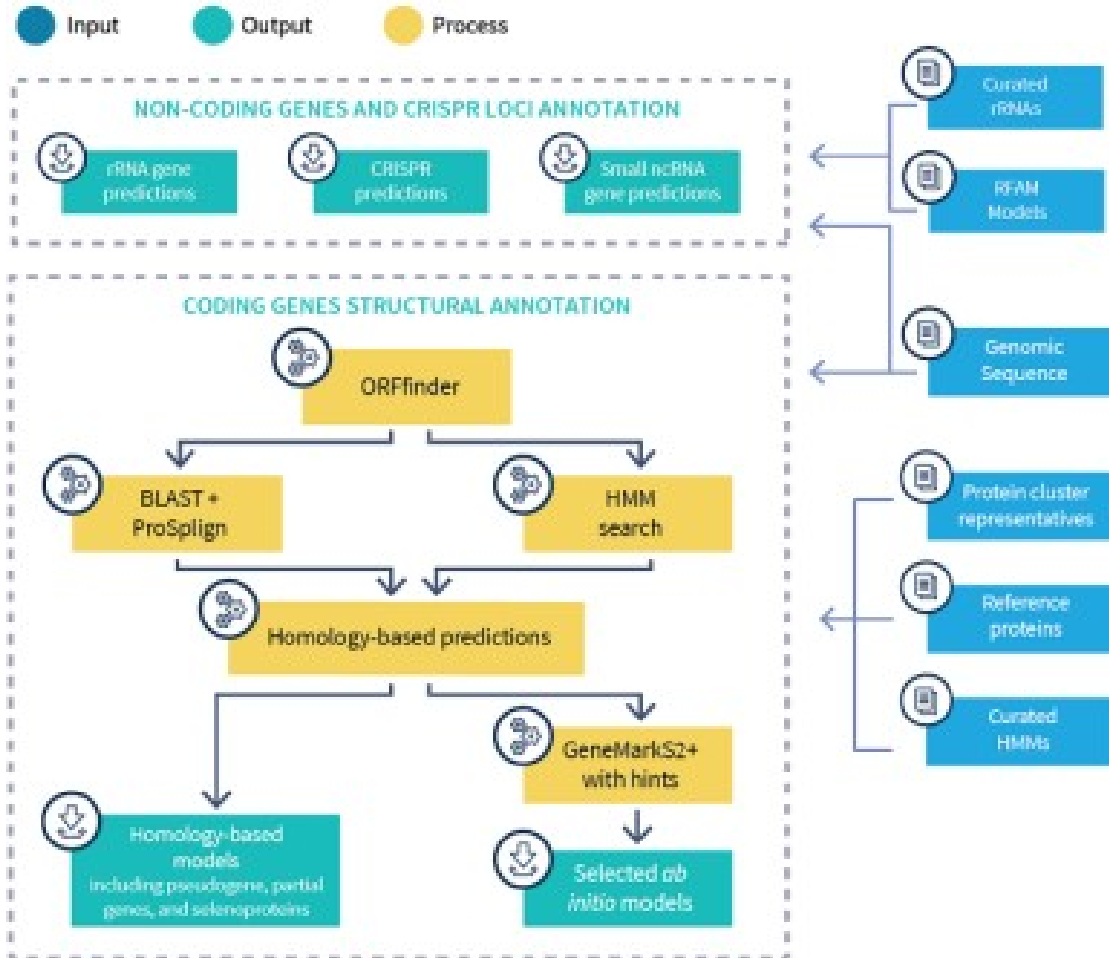- NCBI Prokaryotic Genome Annotation Pipeline

# Prokka

- Illumina BaseSpace app
- Vendor lock

# MicrobeAnnotator

# NCBI Annotation pipeline

# NCBI Annotation pipeline