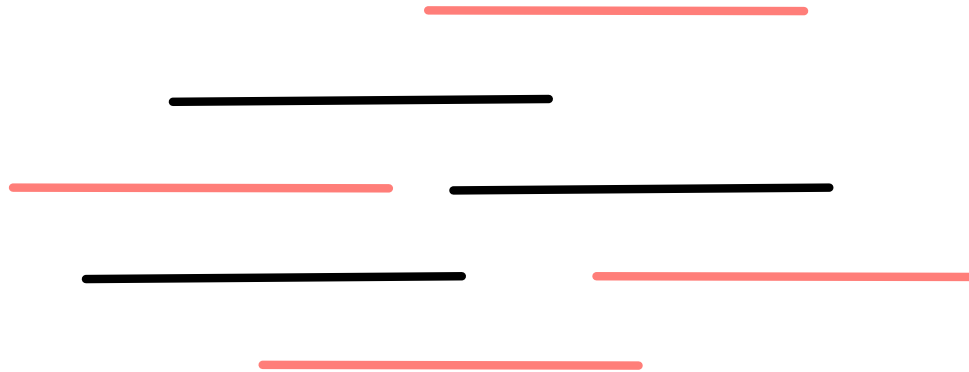


Shotgun sequencing

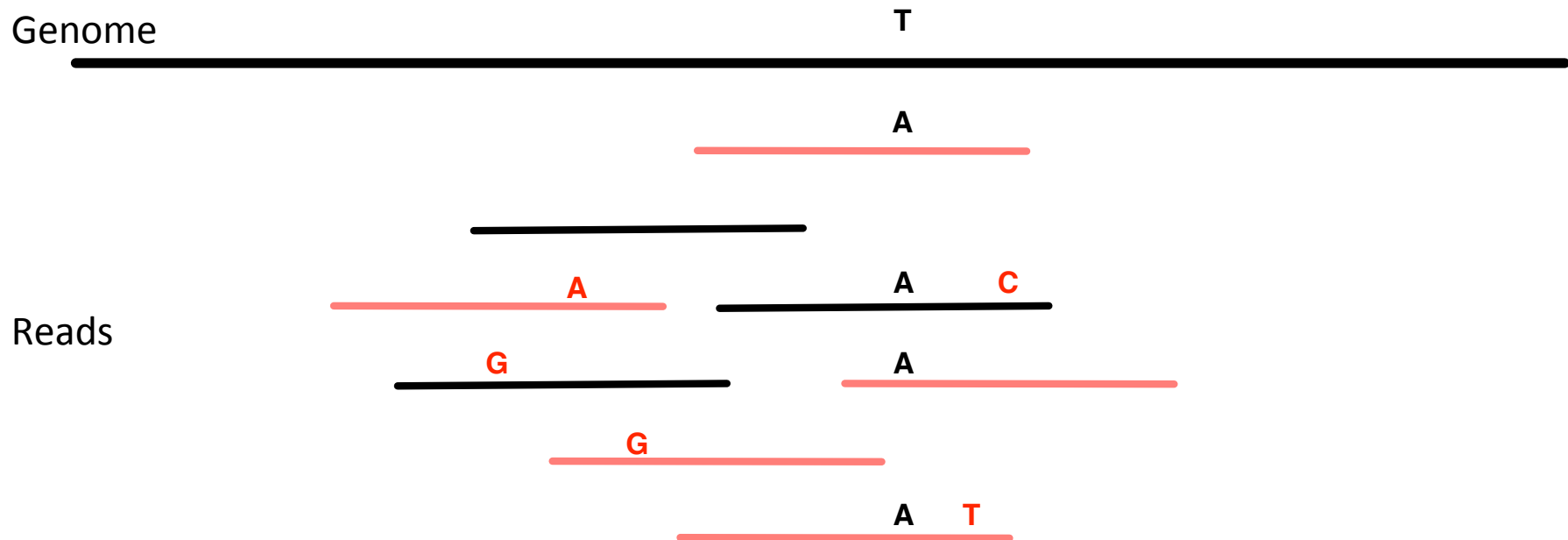
Genome



Reads

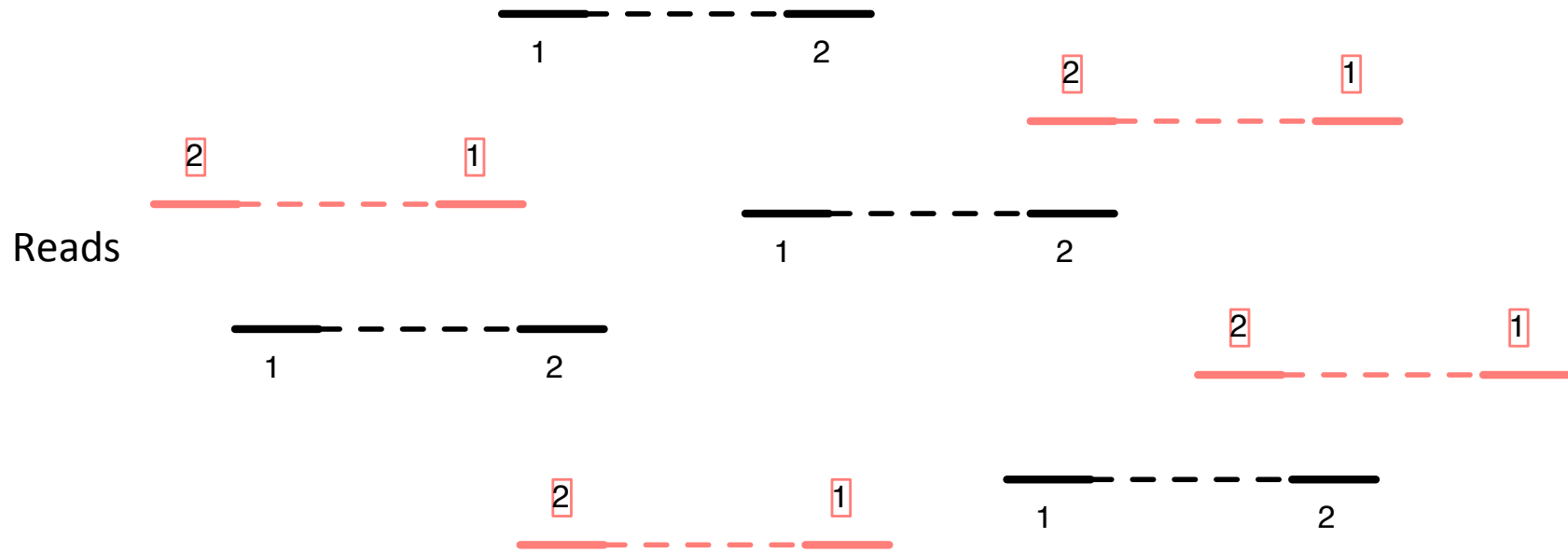


Mismatches: random vs alignable.

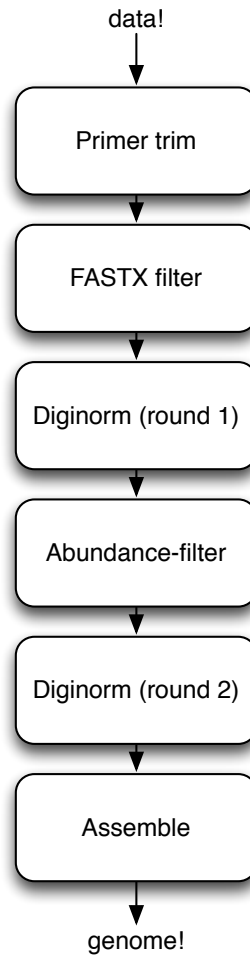


Paired end reads

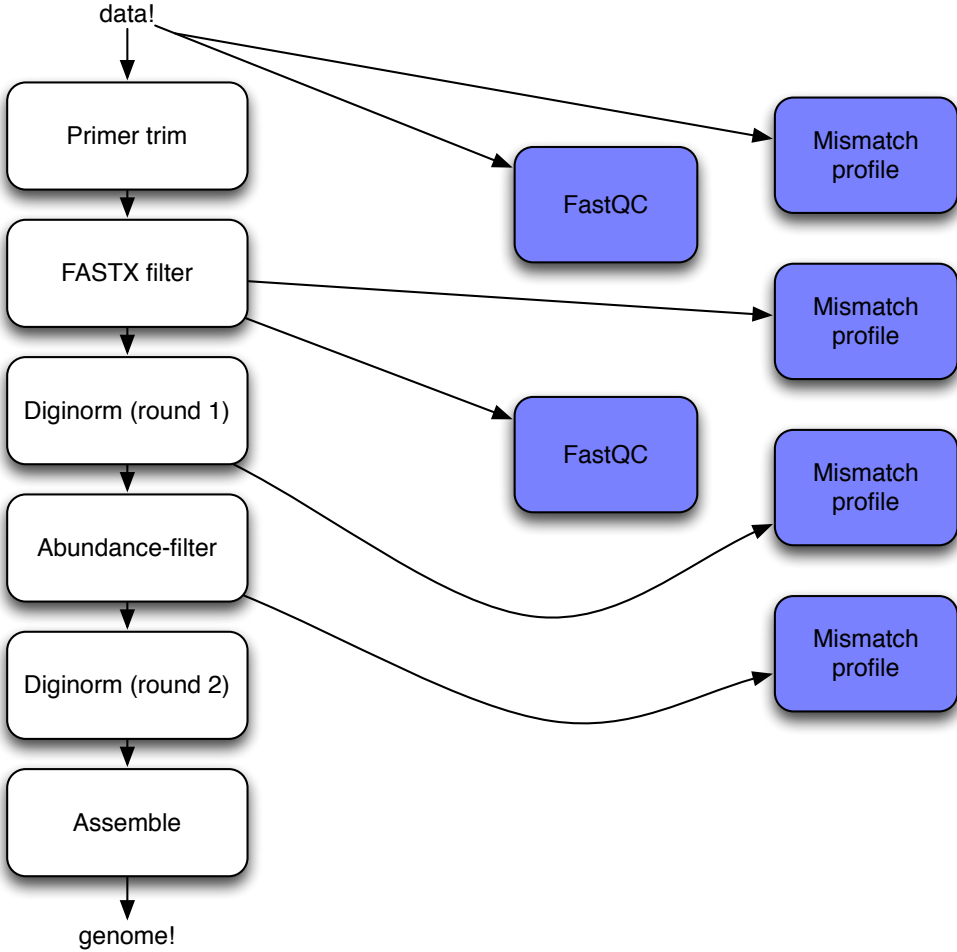
Genome



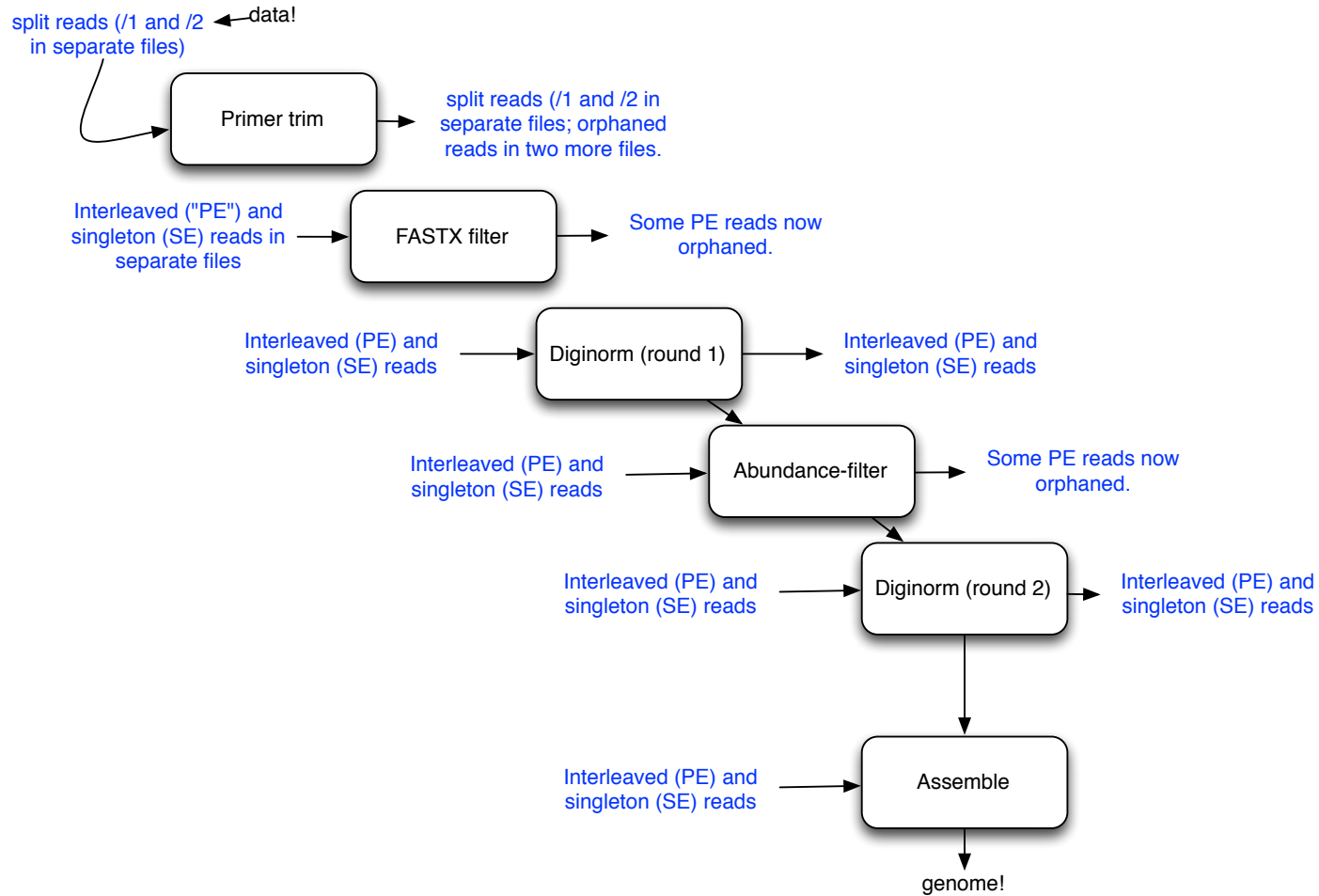
Assembly pipeline for E. coli



Assembly pipeline... with diagnostics



Tracking paired ends!

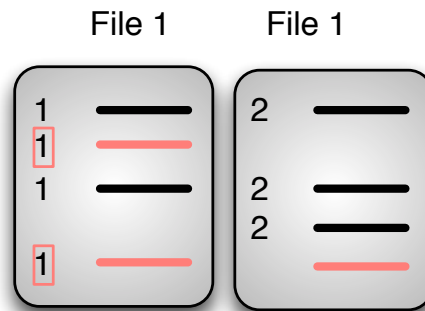


Storing reads: a taxonomy

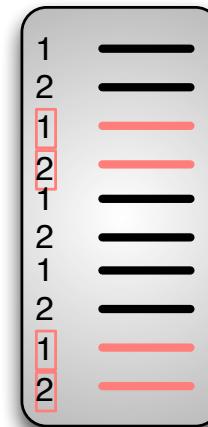
Split, in register



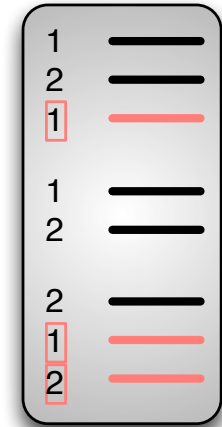
Split w/orphans



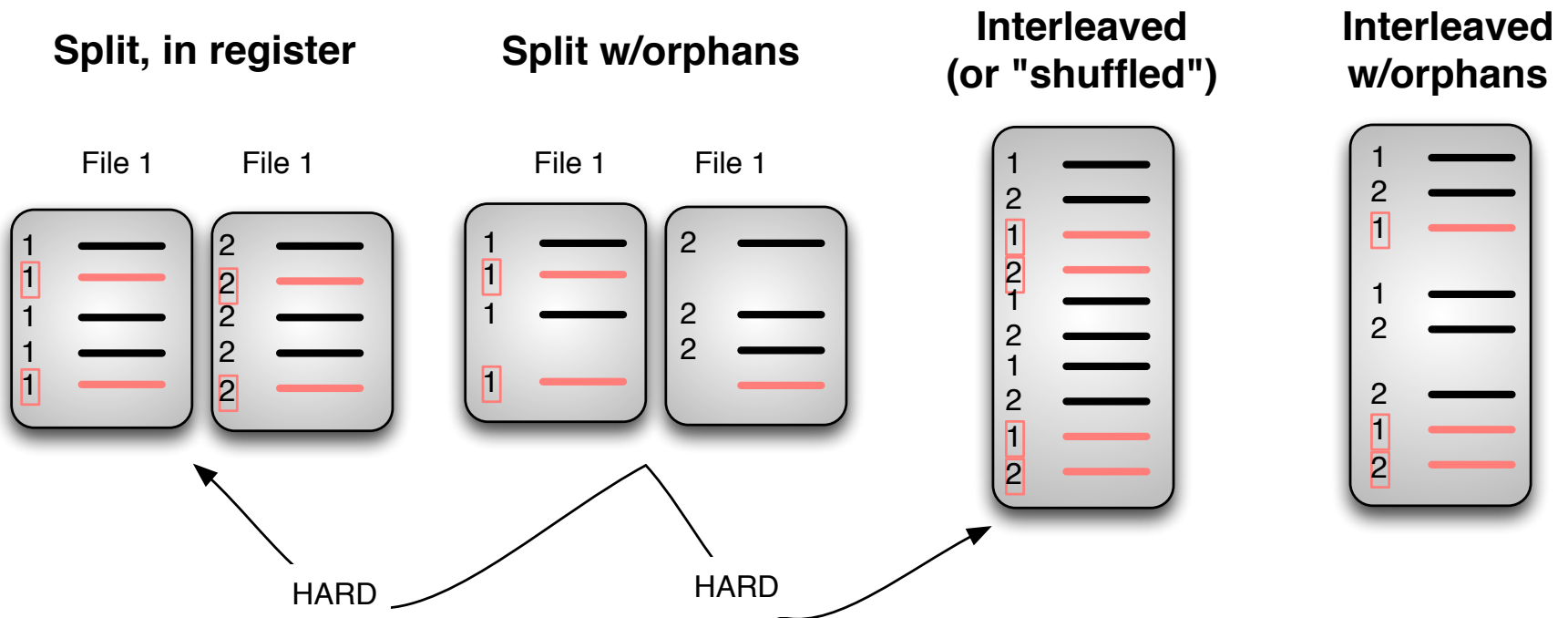
Interleaved
(or "shuffled")



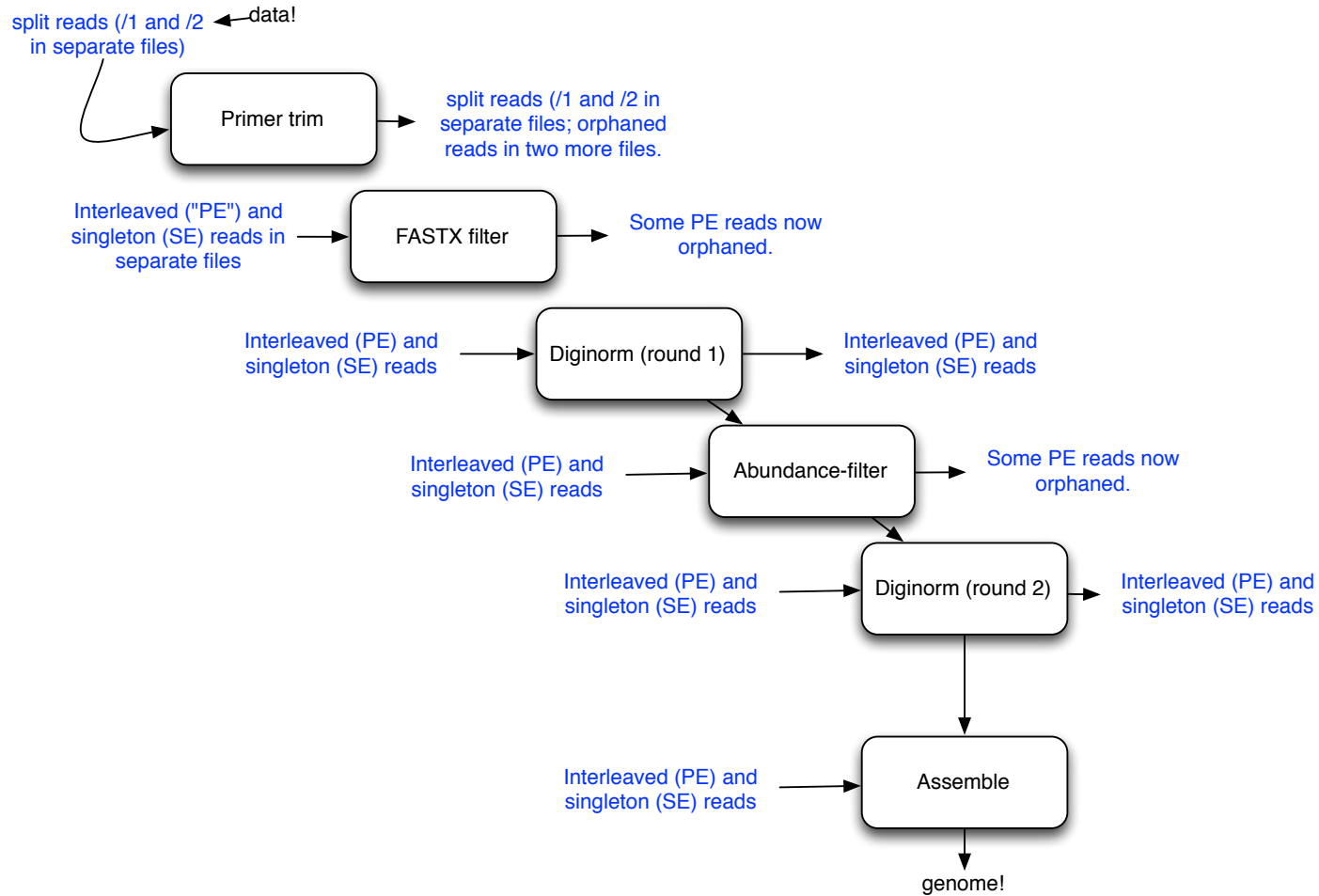
Interleaved
w/orphans



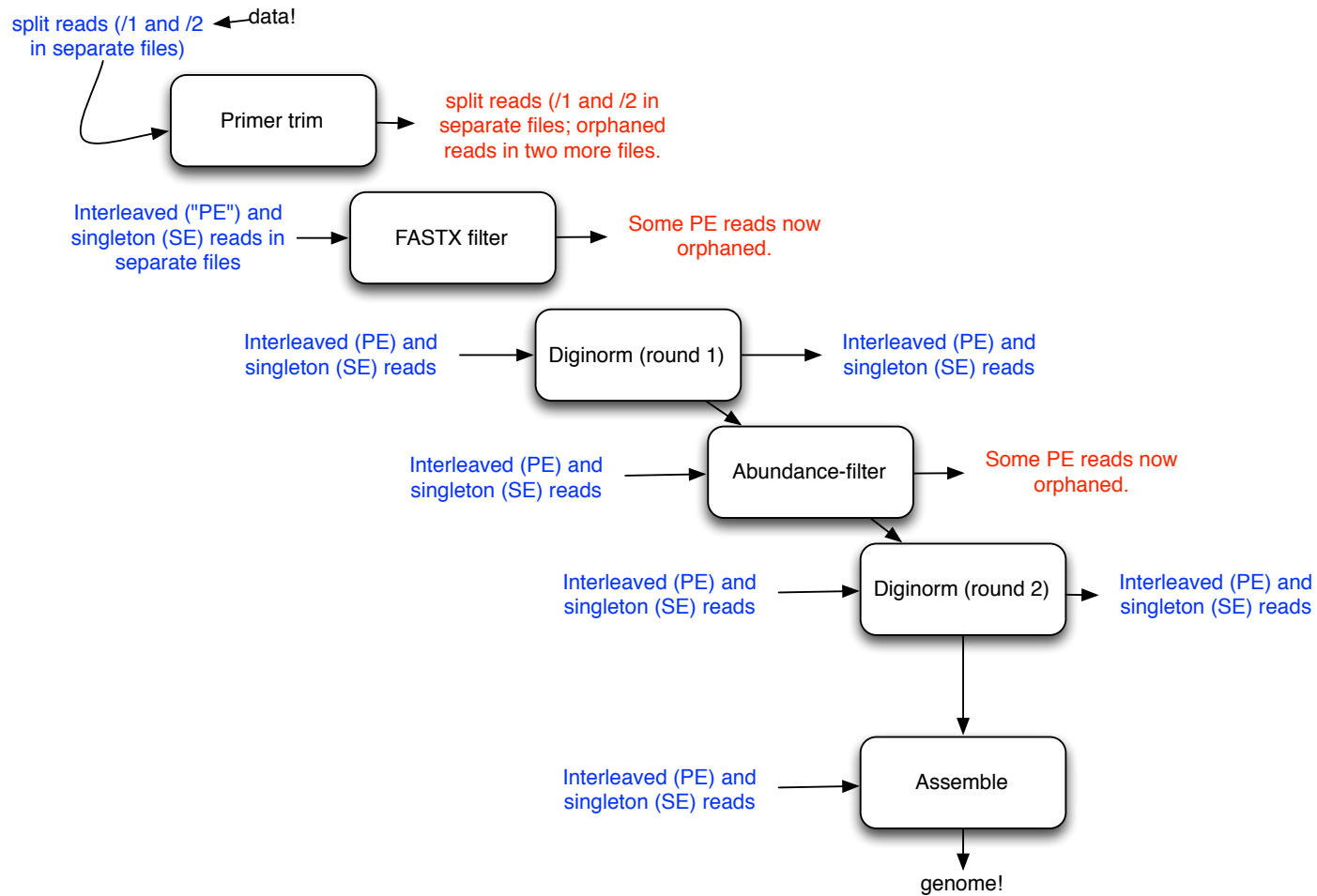
Interconverting can be hard.



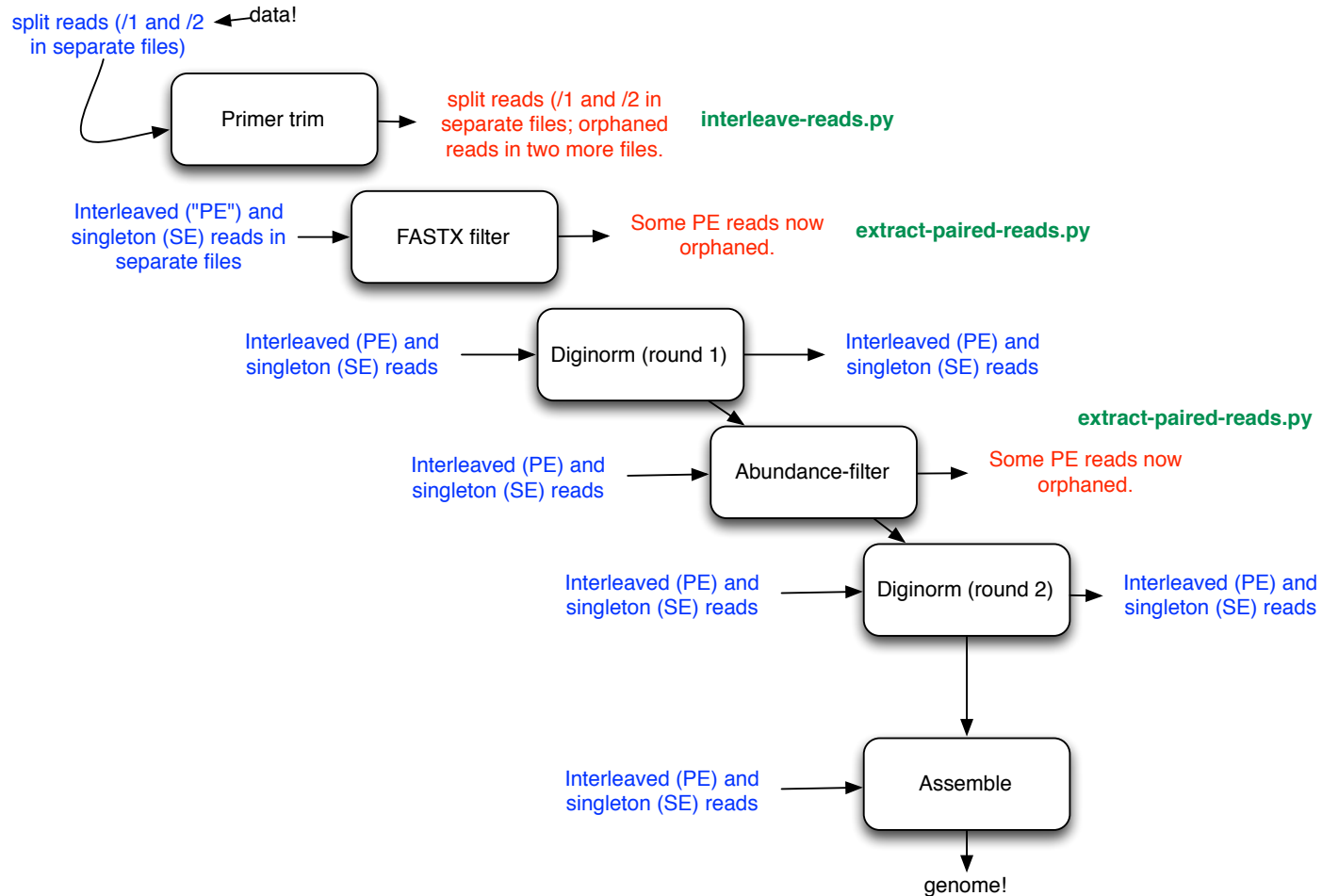
Tracking paired ends:



Must convert.



=>Explicit conversion steps in protocol



Stochastic output

- Velvet (like most assemblers) has stochastic components.
- What this means is that even for same input, same parameters, same machine => different output!?
- Two causes:
 - Random number generator;
 - Multithreading/multiprocessing (multiple “workers”)

Everyone got different results from assembly; did it affect CRP gene?

crp MVLGKPQTDPTLEWFLSHCHIHKYPSTLIHQGEKAETLYYIVKGSVAVLIKDEEGKEM
a MVLGKPQTDPTLEWFLSHCHIHKYPSTLIHQGEKAETLYYIVKGSVAVLIKDEEGKEM
b MVLGKPQTDPTLEWFLSHCHIHKYPSTLIHQGEKAETLYYIVKGSVAVLIKDEEGKEM
d MVLGKPQTDPTLEWFLSHCHIHKYPSTLIHQGEKAETLYYIVKGSVAVLIKDEEGKEM
c MVLGKPQTDPTLEWFLSHCHIHKYPSTLIHQGEKAETLYYIVKGSVAVLIKDEEGKEM

crp ILSYLNQGDFIGELGLFEEGQERSAWVRAKTACEVAEISYKKFRQLIQVNPDILMRLSAQ
a ILSYLNQGDFIGELGLFEEGQERSAWVRAKTACEVAEISYKKFRQLIQVNPDILMRLSAQ
b ILSYLNQGDFIGELGLFEEGQERSAWVRAKTACEVAEISYKKFRQLIQVNPDILMRLSAQ
d ILSYLNQGDFIGELGLFEEGQERSAWVRAKTACEVAEISYKKFRQLIQVNPDILMRLSAQ
c ILSYLNQGDFIGELGLFEEGQERSAWVRAKTACEVAEISYKKFRQLIQVNPDILMRLSAQ

crp MARRLQVTSEKVGNLAFLDVTGRIAQTLNLAQKQPDAMTHPDGMQIKITRQEIGQIVGCS
a MARRLQVTSEKVGNLAFLDVTGRIAQTLNLAQKQPDAMTHPDGMQIKITRQEIGQIVGCS
b MARRLQVTSEKVGNLAFLDVTGRIAQTLNLAQKQPDAMTHPDGMQIKITRQEIGQIVGCS
d MARRLQVTSEKVGNLAFLDVTGRIAQTLNLAQKQPDAMTHPDGMQIKITRQEIGQIVGCS
c MARRLQVTSEKVGNLAFLDVTGRIAQTLNLAQKQPDAMTHPDGMQIKITRQEIGQIVGCS

crp RETVGRILKMLEDQNLISAHGKTIVVYGTR
a RETVGRILKMLEDQNLISAHGKTIVVYGTR
b RETVGRILKMLEDQNLISAHGKTIVVYGTR
d RETVGRI-----
c RETVGRILKMLEDQNLISAHGKTIVVYGTR

Program output in large-scale analyses

- You can almost always get something that makes *some* sense, i.e. isn't wrong.
- But: Insensitive? Incomplete? Biased?

Next two weeks in bioinfo

- Annotating genome
- Mapping & variant visualization
- Mapping & Quantification