

# Creating supermatrices and phylogenetic datasets with **SuperCRUNCH**

**Daniel Portik**

Postdoctoral Researcher  
University of Arizona

# Phylogenetic Data Workflows



Handful of tools available for assembling phylogenetic data

# Phylogenetic Data Workflows

Handful of tools available for assembling phylogenetic data

Each has great features!

<b>PhyLoTa</b>	(Sanderson et al. 2008)
<b>PHLAWD</b>	(Smith et al. 2009)
<b>phyloGenerator</b>	(Pearse & Purvis 2013)
<b>SUMAC</b>	(Freyman, 2015)
<b>SUPERSMART</b>	(Antonelli et al. 2017)
<b>PhylotaR</b>	(Bennett et al. 2018)
<b>PyPHLAWD</b>	(Smith & Walker 2018)

# Phylogenetic Data Workflows

## Common limitations for my work needs

- Sequences fetched using GenBank release database
- Reliance on NCBI Taxonomy for identifying sequences
- Automated clustering of whole sequence sets
- Black-box methods for sequence filtering and selection

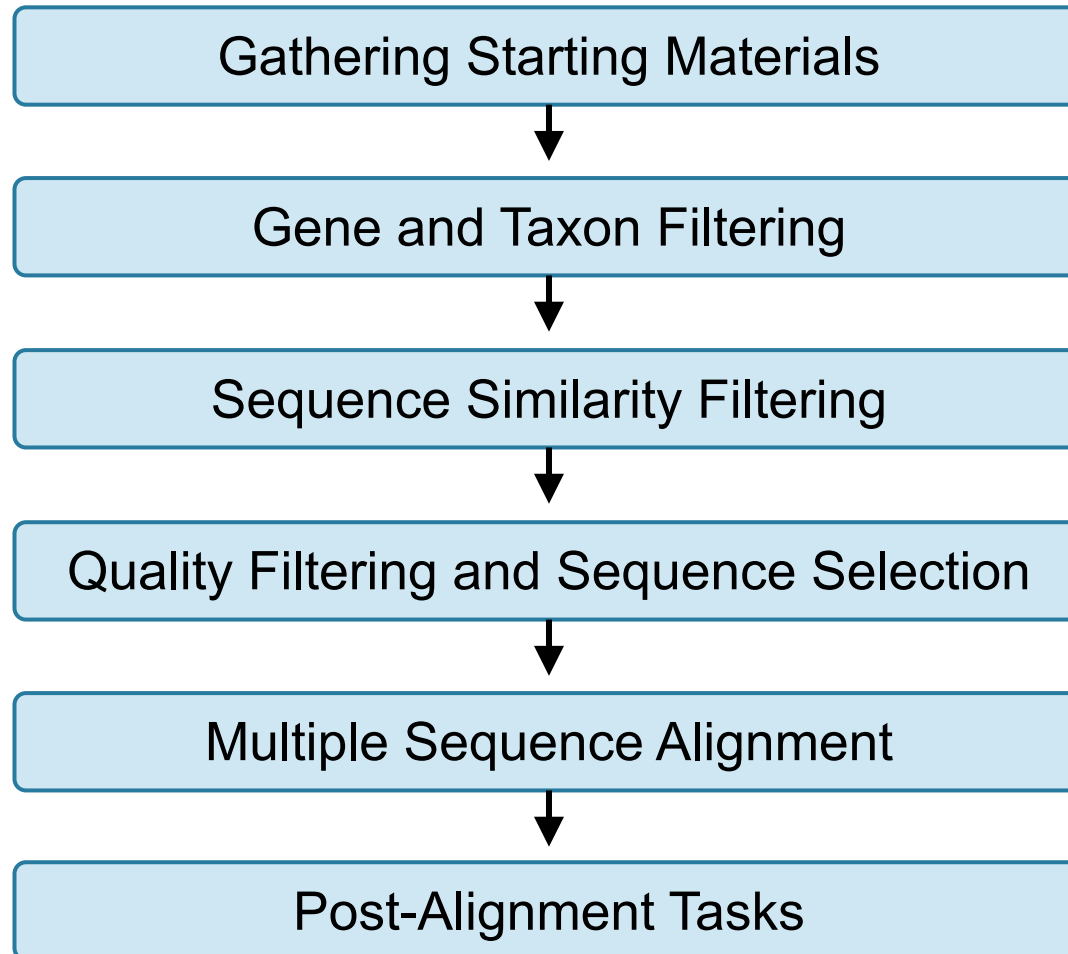
# SuperCRUNCH

FOR PHYLOGENETIC DATA

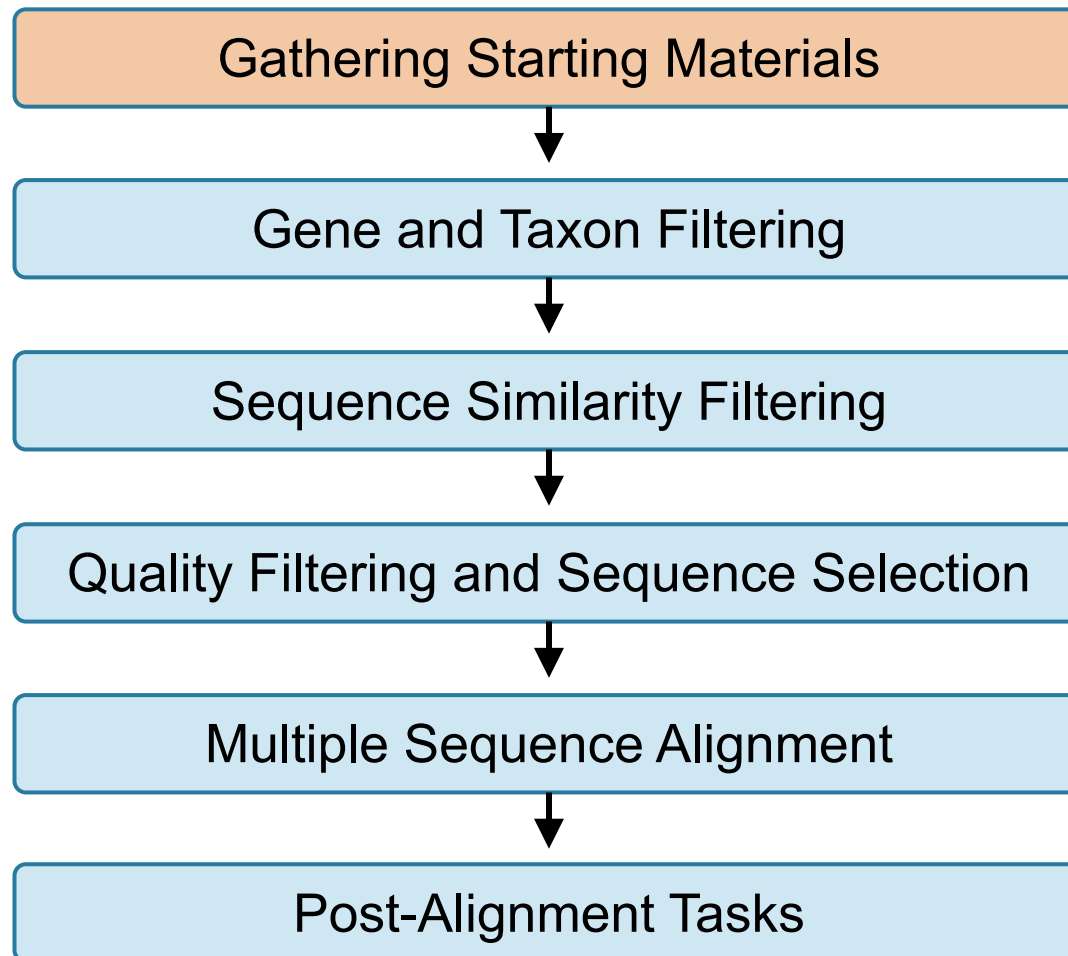
The logo graphic consists of seven horizontal yellow bars of varying lengths, stacked vertically to the right of the text.

- Bioinformatic toolkit for phylogenetic data
- Allows ANY fasta-formatted sequence data to be used
- Targeted searches based on list of taxa and genes
- Transparent methods for orthology-filtering and sequence selection, including detailed output files
- Can be used to assemble traditional supermatrices as well as population-level datasets
- Can process large sequence capture datasets (UCE's)

# SuperCRUNCH Workflow



# SuperCRUNCH Workflow



# Gathering Starting Materials

Sequence Data

Taxon List

Gene List



# Gathering Starting Materials

Sequence Data

Taxon List

Gene List

NCBI/GenBank  
sequences

and/or

Locally generated  
sequences

# Gathering Starting Materials

Sequence Data

Taxon List

Gene List

NCBI/GenBank  
sequences

and/or

Locally generated  
sequences

Fasta format:

```
>KP820543.1 Callisaurus draconoides voucher R45 aryl hydrocarbon  
receptor (anr) gene, partial cds  
CACAAATGAGAAAGCCTTGATAAACCGTGATCGGACTTTGCCACTCGTTGAAGAAATAGATGAGAGCT...
```

```
>KP820544.1 Urosaurus ornatus voucher UWBM7587 aryl hydrocarbon  
receptor (anr) gene, partial cds  
TAAATCTCCTTTGAAAGGAACCTTTTTGTGGACACCAGGGATGAATTAGGTAATGTAATGGCCAGA...
```

# Gathering Starting Materials

Sequence Data

Taxon List

Gene List

NCBI/GenBank  
sequences

and/or

Locally generated  
sequences

Fasta format:

> **Unique Identifier**   **Taxon Label**   **Description**

```
>KP820543.1 Callisaurus draconoides voucher R45 aryl hydrocarbon  
receptor (anr) gene, partial cds  
CACAAATGAGAAAGCCTTGATAAACCGTGATCGGACTTTGCCACTCGTTGAAGAAATAGATGAGAGCT...
```

```
>KP820544.1 Urosaurus ornatus voucher UWM7587 aryl hydrocarbon  
receptor (anr) gene, partial cds  
TAAATCTCCTTTGAAAGGAACCTTTTTGTGGACACCAGGGATGAATTAGGTAATGTAATGGCCAGA...
```

# Gathering Starting Materials

Sequence Data

Taxon List

Gene List

Text file with species  
and subspecies labels

```
Callisaurus draconoides crinitus  
Callisaurus draconoides draconoides  
Callisaurus draconoides inusitanus  
Cophosaurus texanus  
Holbrookia maculata  
Uma exsul  
Uma inornata  
Uma notata  
Uma paraphygas  
Uma rufopunctata  
Uma scoparia
```

# Gathering Starting Materials

Sequence Data

Taxon List

Gene List

Text file with species  
and subspecies labels

List Sources

```
Callisaurus draconoides crinitus  
Callisaurus draconoides draconoides  
Callisaurus draconoides inusitanus  
Cophosaurus texanus  
Holbrookia maculata  
Uma exsul  
Uma inornata  
Uma notata  
Uma paraphygas  
Uma rufopunctata  
Uma scoparia
```

- Taxonomic Databases
- NCBI Taxonomy Browser
- Directly from sequences  
(via SuperCRUNCH)

# Gathering Starting Materials

Sequence Data

Taxon List

Gene List

Tab-delimited text file  
with search terms

- Search terms include a gene abbreviation and description
- Can include any number and combination of genes (mtDNA, nuclear, seq-cap)
- Tetrapod 5k UCE set and other example files on github

# Gathering Starting Materials

Sequence Data

Taxon List

Gene List

Tab-delimited text file  
with search terms

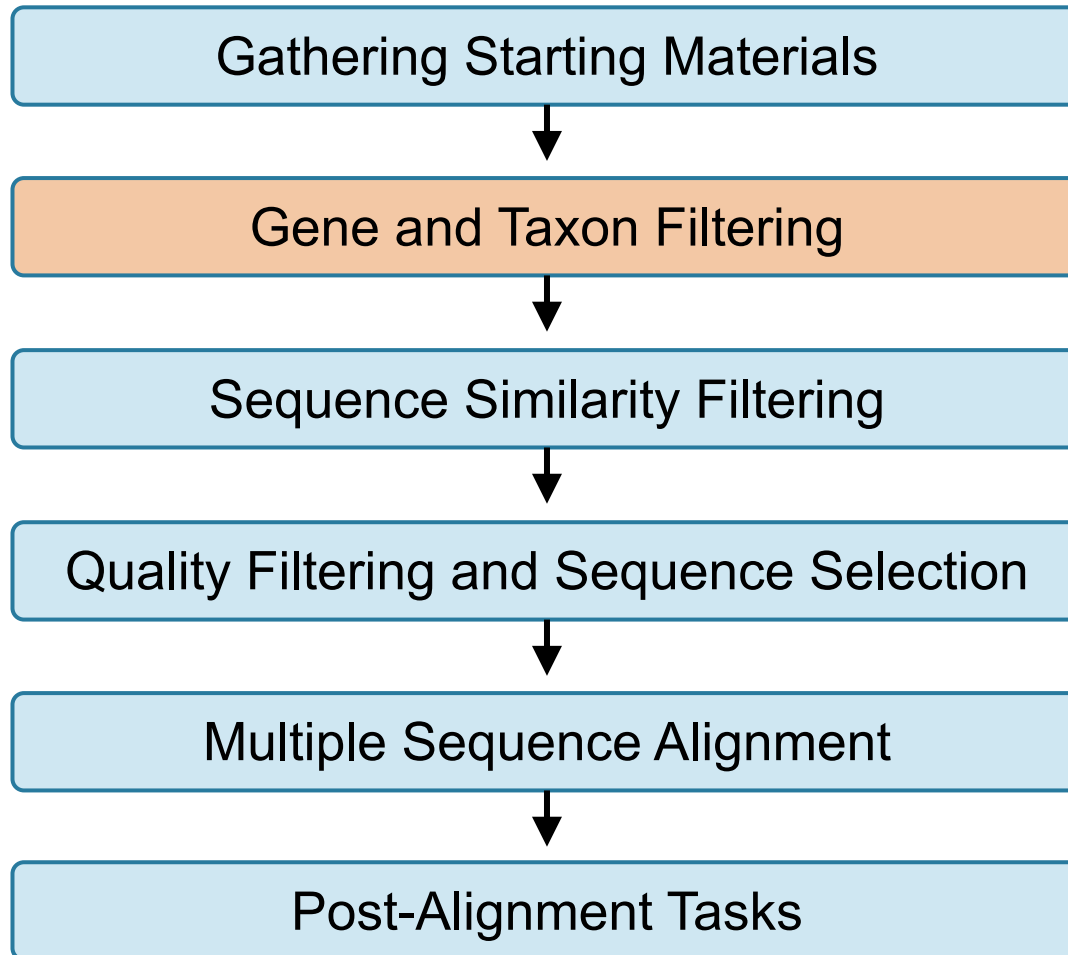
Label

Abbreviation

Description

CMOS	CMOS;C-MOS	oocyte maturation factor
CXCR4	CXCR4	chemokine C-X-C motif receptor 4;C-X-C ...
DLL1	DLL1;DLL	distal-less
DNAH3	DNAH3	dynein axonemal heavy chain 3;dynein ax...
ECEL1	ECEL1;ECEL	endothelin converting;endothelin conver...
ENC1	ENC1	ectodermal neural cortex 1
EXPH5	EXPH5	exophilin;exophilin 5;exophilin-5;exoph...
FSHR	FSHR	follicle stimulating hormone;follicle-s...

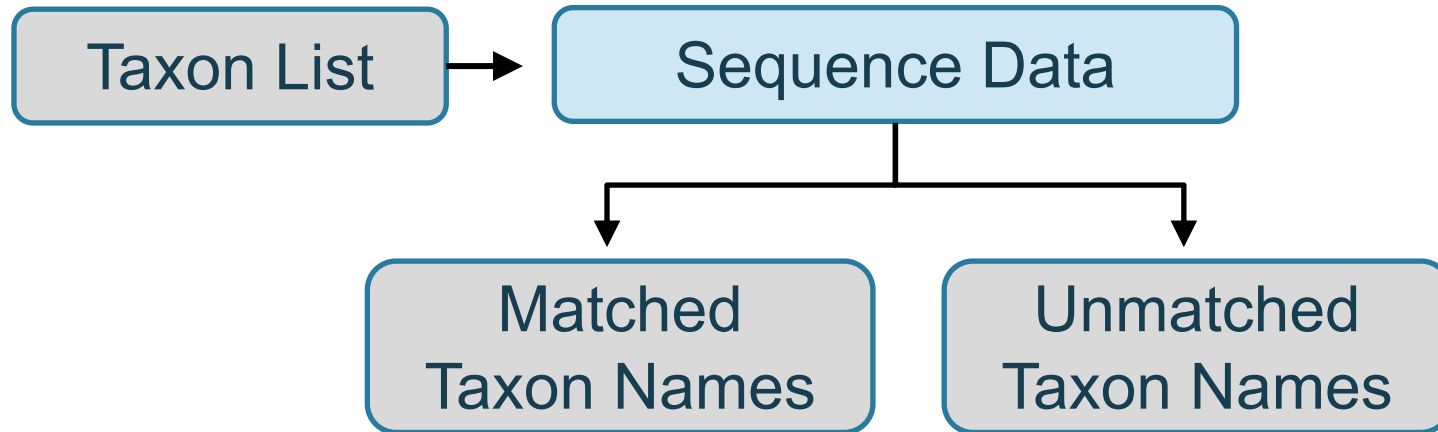
# SuperCRUNCH Workflow





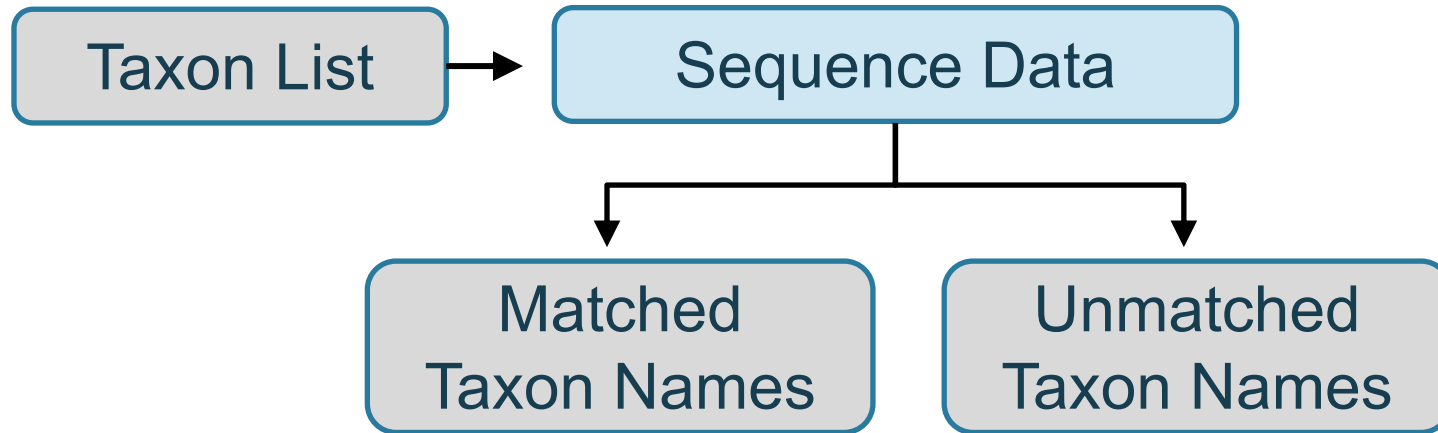
# Gene and Taxon Filtering

## Optional assessment and cleaning step



# Gene and Taxon Filtering

## Optional assessment and cleaning step

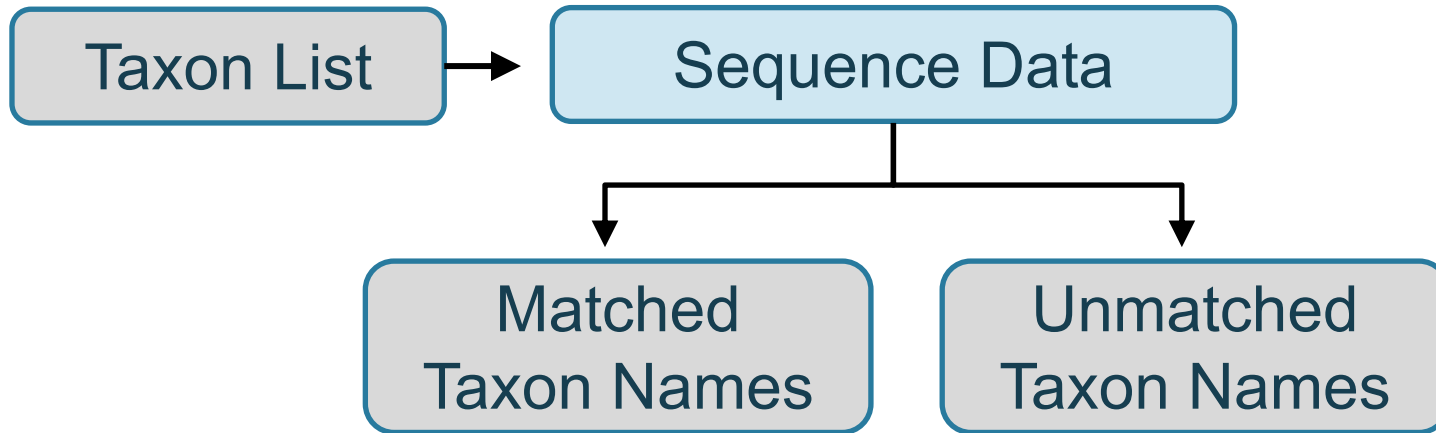


## Bad Records

A.porcus mitochondrial  
Calumma sp.  
Jp 2016507218-a/5  
Unverified calotes

# Gene and Taxon Filtering

## Optional assessment and cleaning step



Good records, but need to be corrected

Uromastix acanthinura



Uromastyx acanthinura

Causus defilippi



Causus defilippii

Vipera raddei



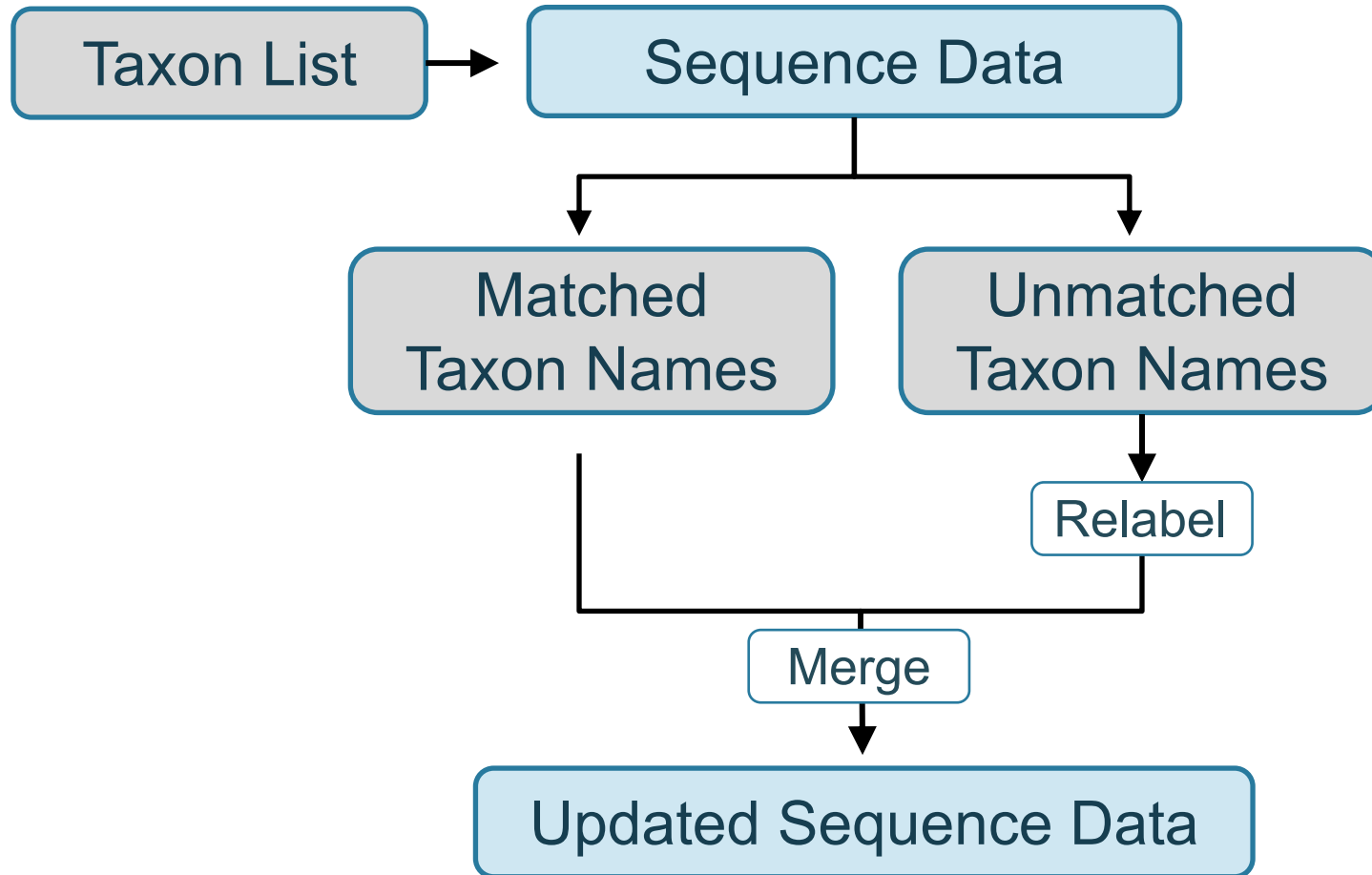
Montivipera raddei

Spelling errors

Change in taxonomy

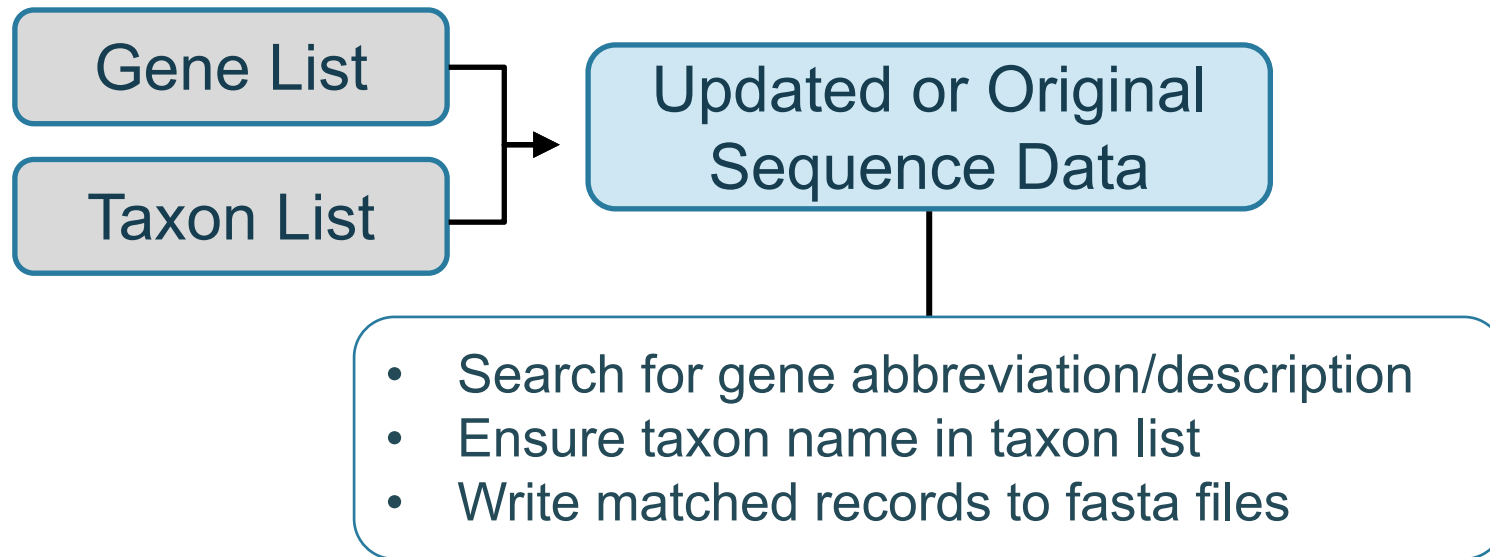
# Gene and Taxon Filtering

## Optional assessment and cleaning step



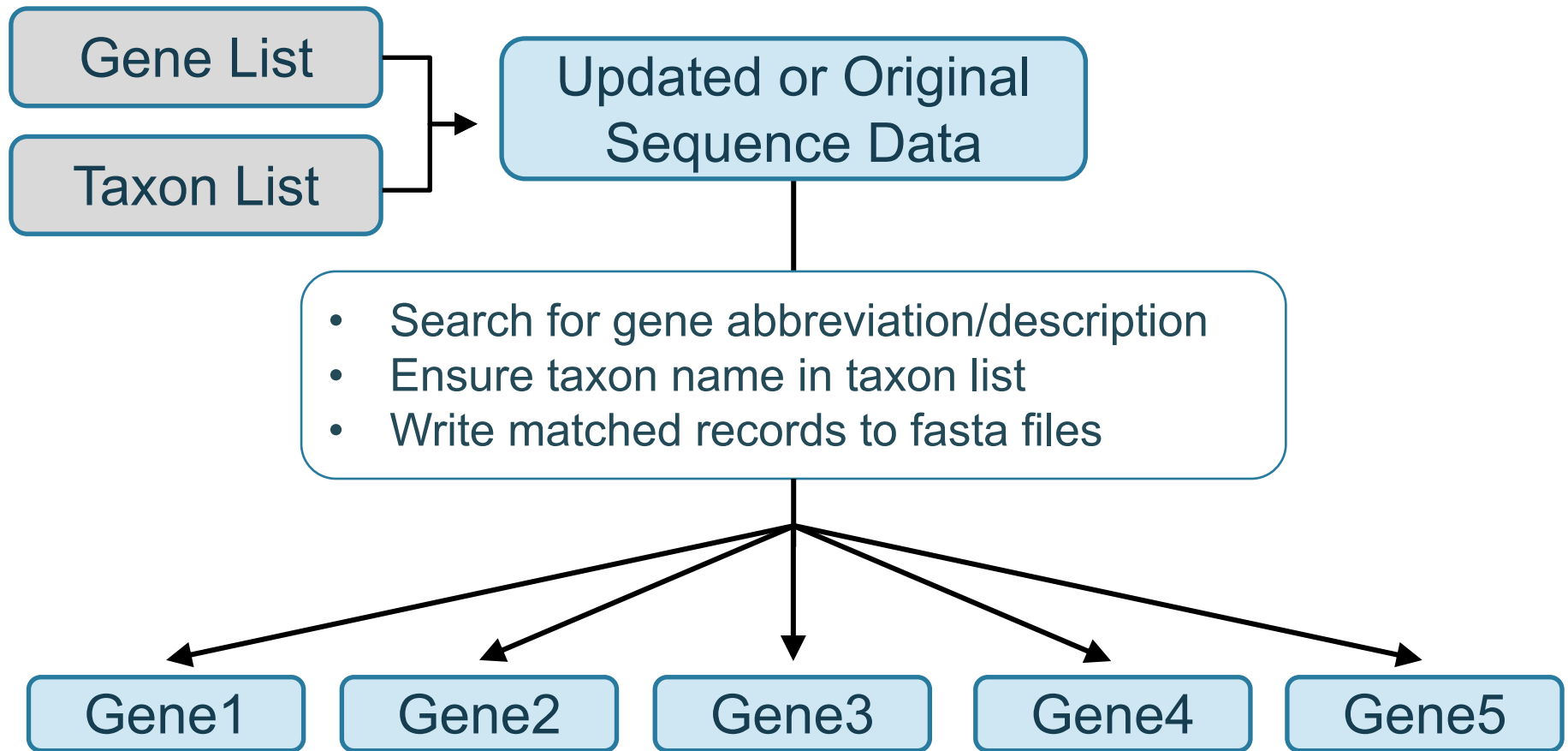
# Gene and Taxon Filtering

## Parsing Loci

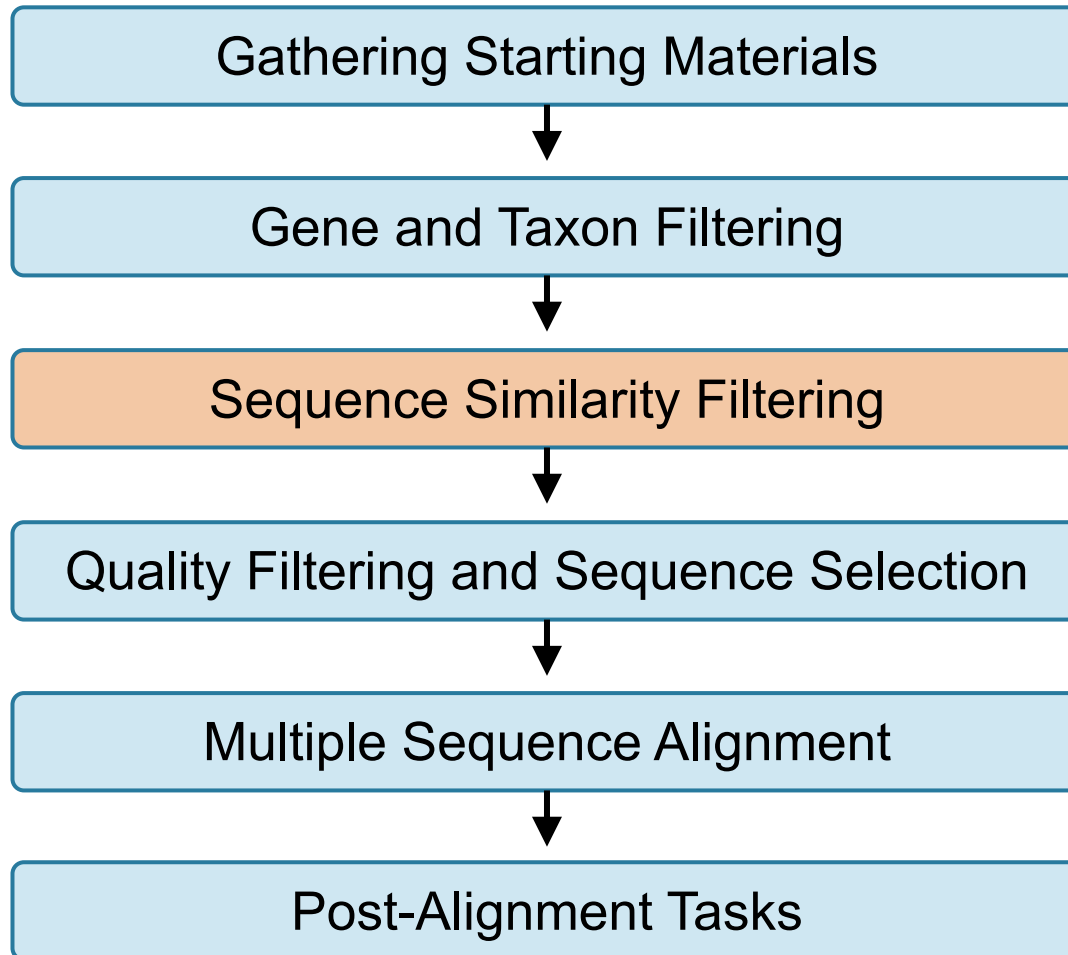


# Gene and Taxon Filtering

## Parsing Loci



# SuperCRUNCH Workflow



# Sequence Similarity Filtering

---

**Use BLASTn to identify and extract target sequence**



# Sequence Similarity Filtering

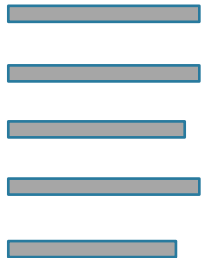
**Use BLASTn to identify and extract target sequence**

Input  
Sequence

Reference  
Sequences



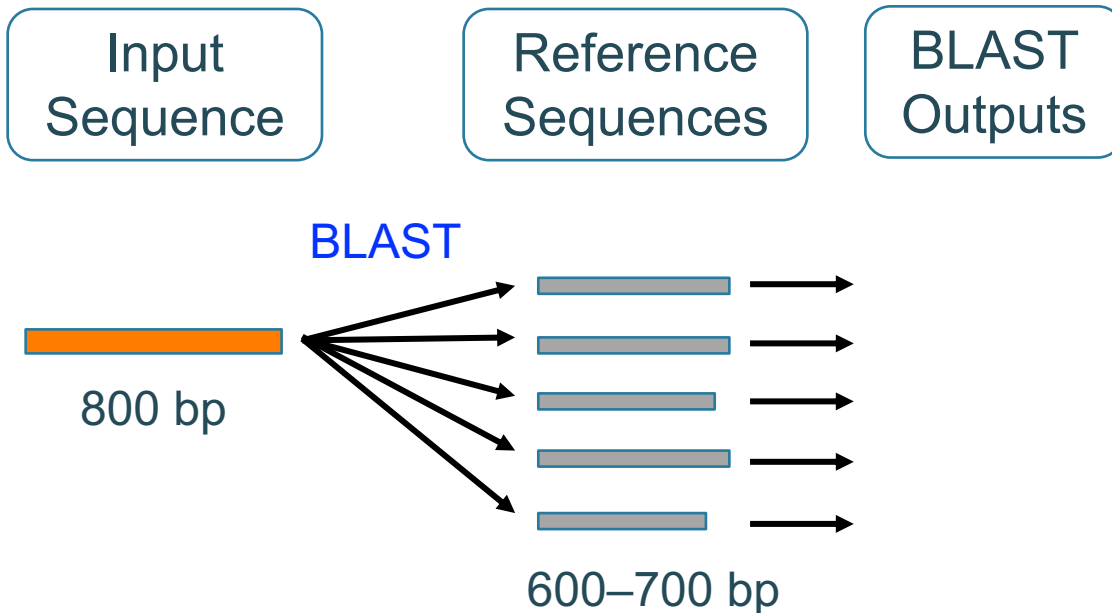
800 bp



600–700 bp

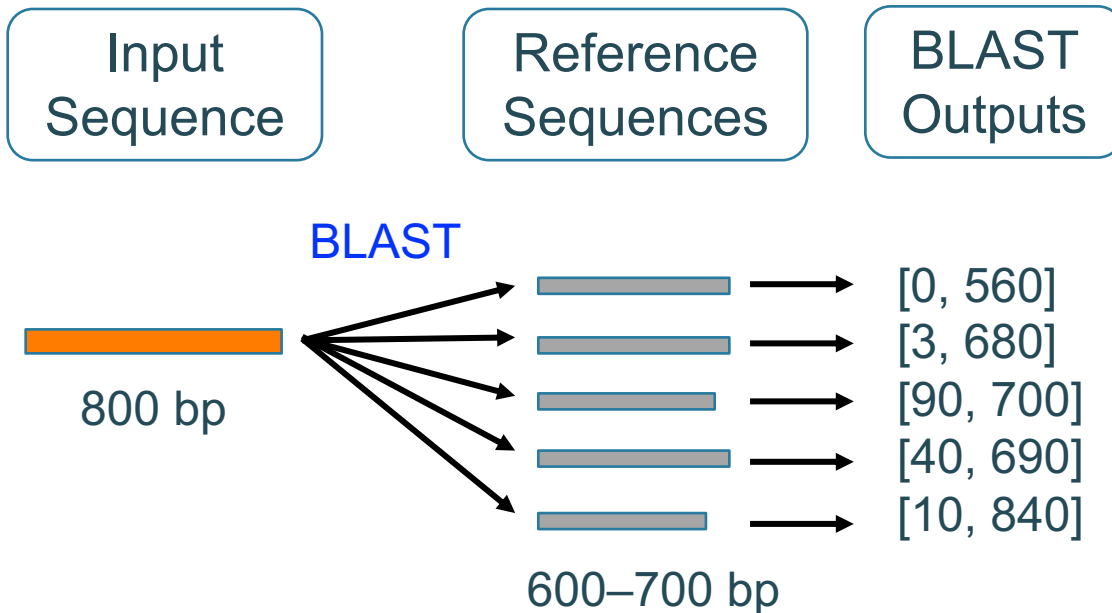
# Sequence Similarity Filtering

**Use BLASTn to identify and extract target sequence**



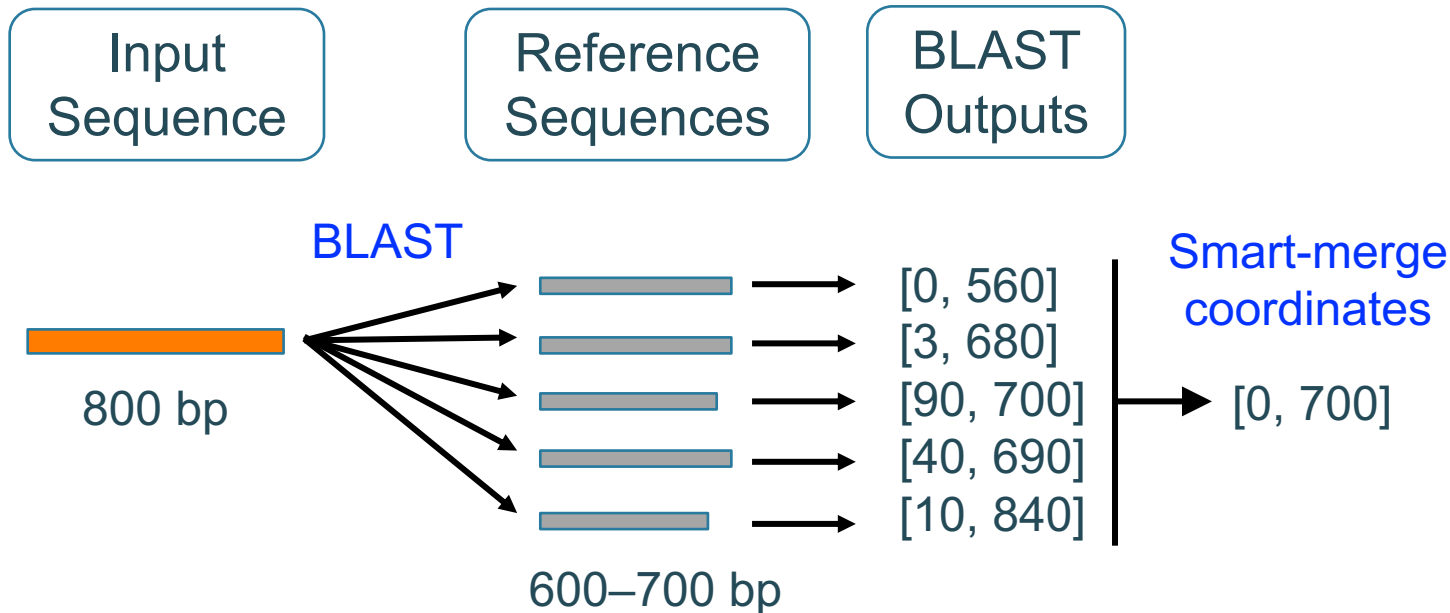
# Sequence Similarity Filtering

**Use BLASTn to identify and extract target sequence**



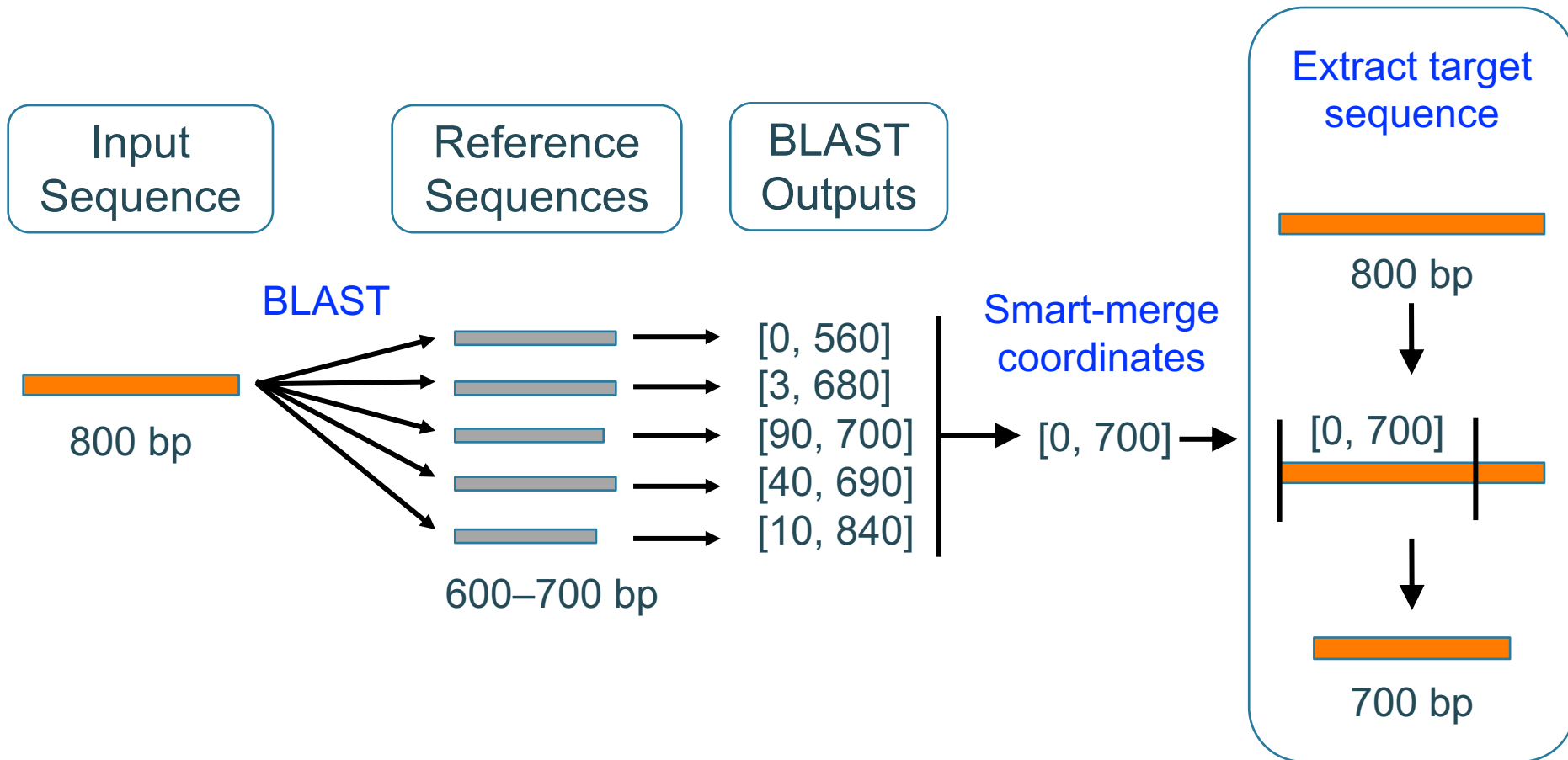
# Sequence Similarity Filtering

**Use BLASTn to identify and extract target sequence**



# Sequence Similarity Filtering

**Use BLASTn to identify and extract target sequence**



# Sequence Similarity Filtering

**Use BLASTn to identify and extract target sequence**

'Simple' Records

'Complex' Records

# Sequence Similarity Filtering

**Use BLASTn to identify and extract target sequence**

‘Simple’ Records



1,200 bp

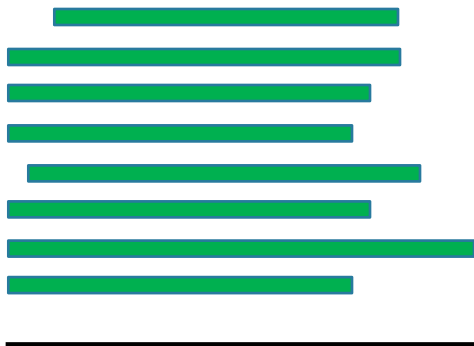
‘Complex’ Records

- Single gene
- Same region obtained
- Minor length variation

# Sequence Similarity Filtering

Use BLASTn to identify and extract target sequence

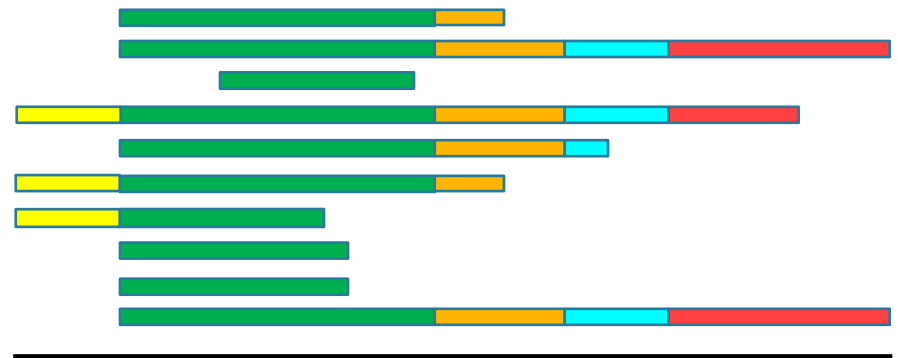
## 'Simple' Records



1,200 bp

- Single gene
- Same region obtained
- Minor length variation

## 'Complex' Records



6,000 bp

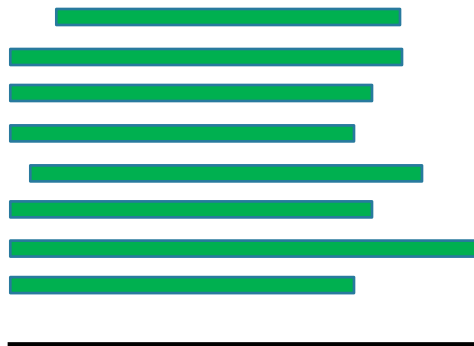
- Often multiple genes
- Variation in region obtained
- Major length variation



# Sequence Similarity Filtering

**Use BLASTn to identify and extract target sequence**

## 'Simple' Records

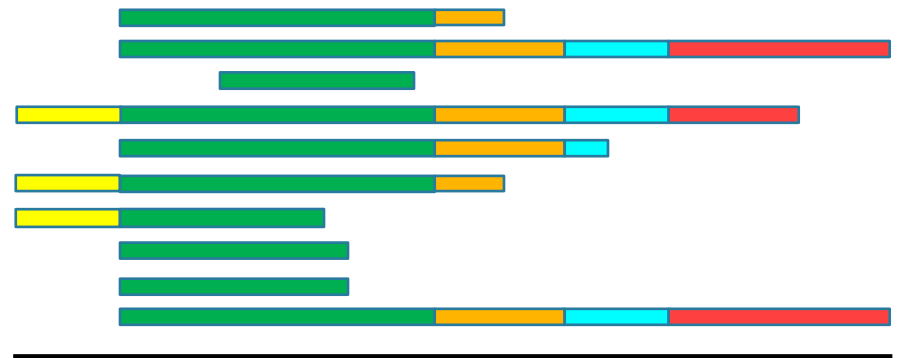


1,200 bp

- Single gene
- Same region obtained
- Minor length variation

**Automatic selection of  
reference sequences**

## 'Complex' Records



6,000 bp

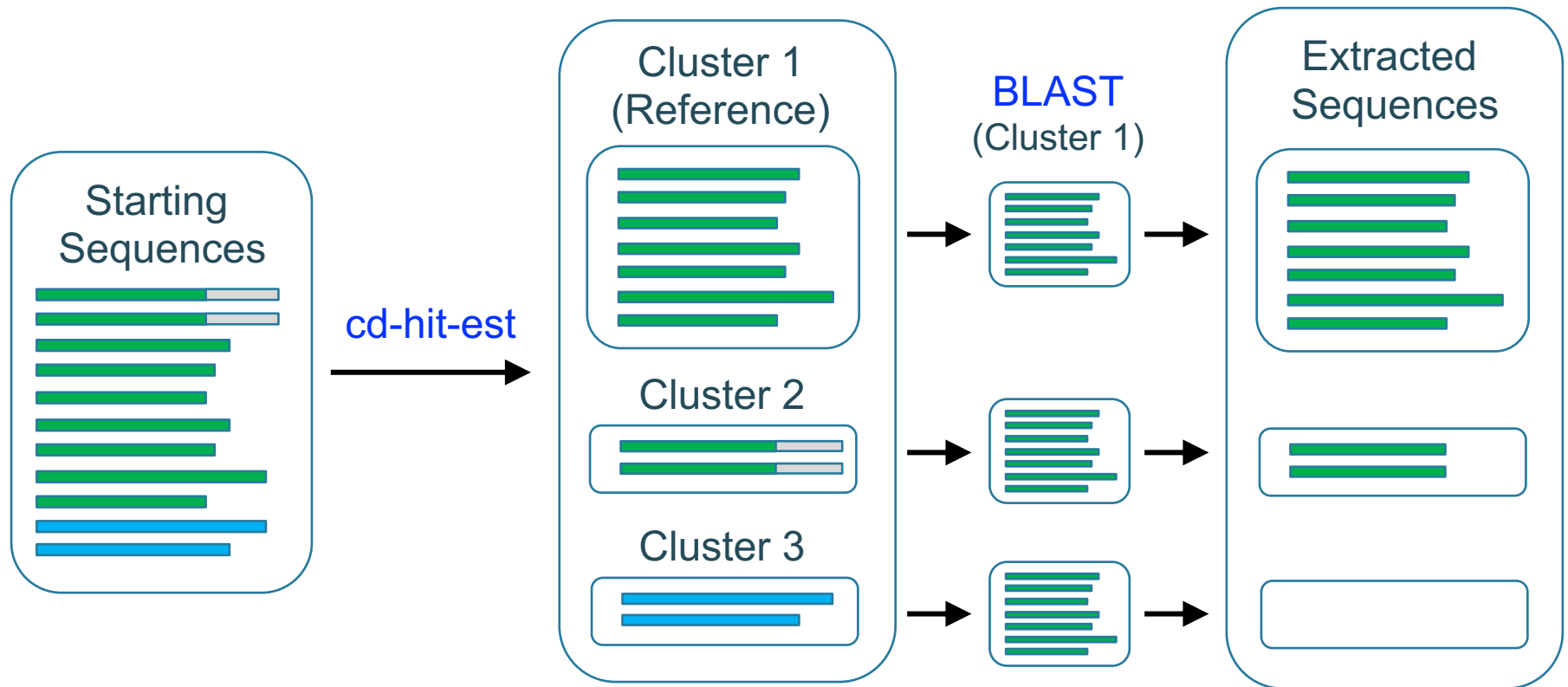
- Often multiple genes
- Variation in region obtained
- Major length variation

**Requires a user-supplied  
reference sequences**

# Sequence Similarity Filtering

'Simple' Records

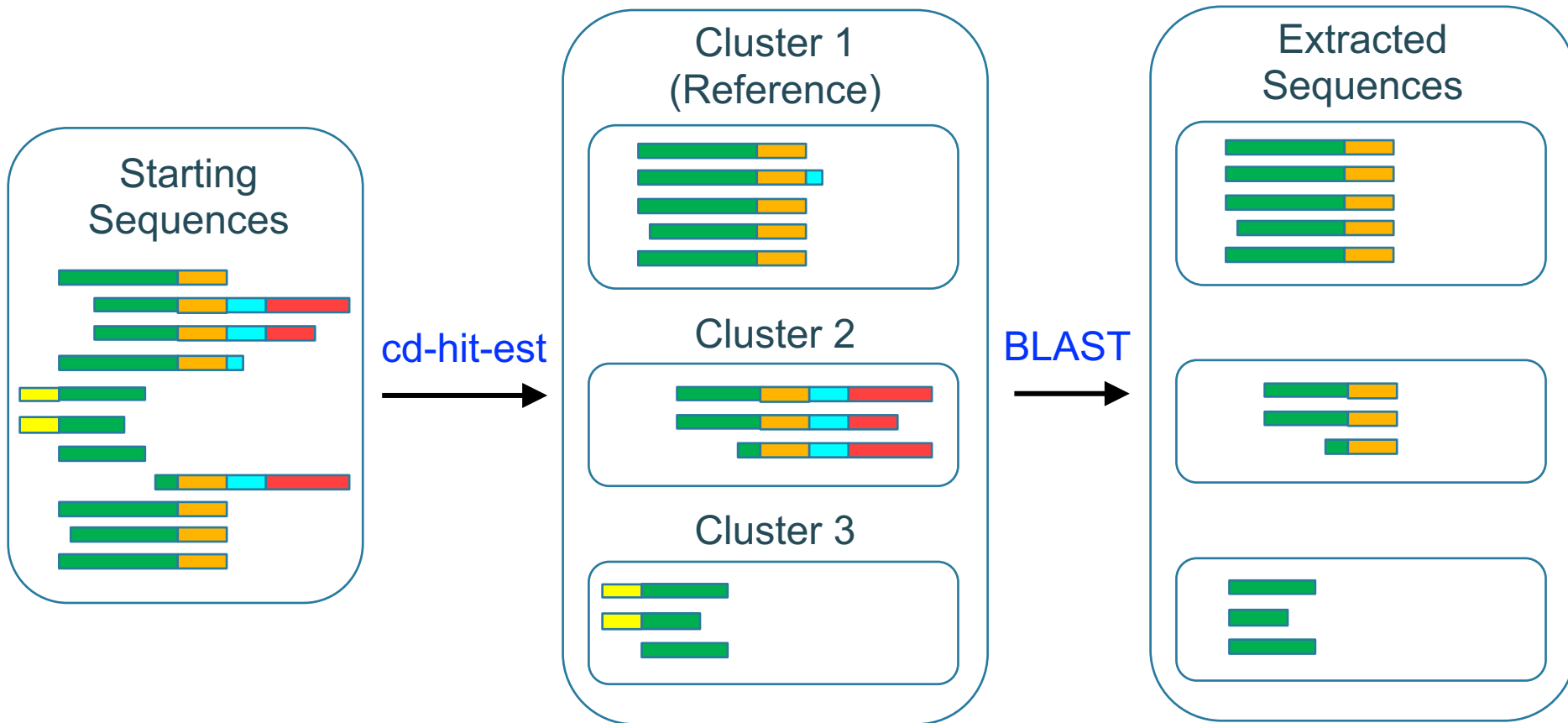
**Automatic selection of reference sequences**



# Sequence Similarity Filtering

**'Complex' Records**

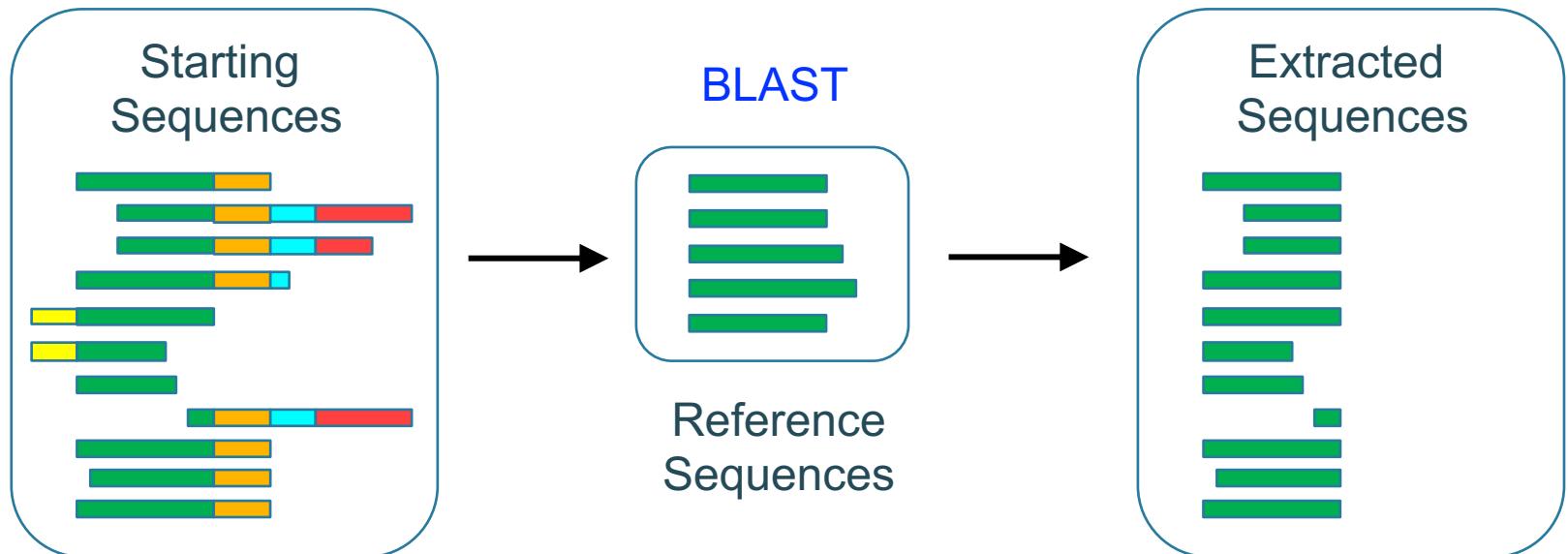
**Automatic selection of reference sequences**



# Sequence Similarity Filtering

**'Complex' Records**

**Supply validated reference sequences**



# Sequence Similarity Filtering

- Ability to specify BLAST algorithm for searches
  - blastn, megablast, dc-megablast

# Sequence Similarity Filtering

- Ability to specify BLAST algorithm for searches
  - blastn, megablast, dc-megablast
- Detailed outputs provide complete transparency
  - sequences discarded, info for extracted sequences

# Sequence Similarity Filtering

- Ability to specify BLAST algorithm for searches
  - blastn, megablast, dc-megablast
- Detailed outputs provide complete transparency
  - sequences discarded, info for extracted sequences

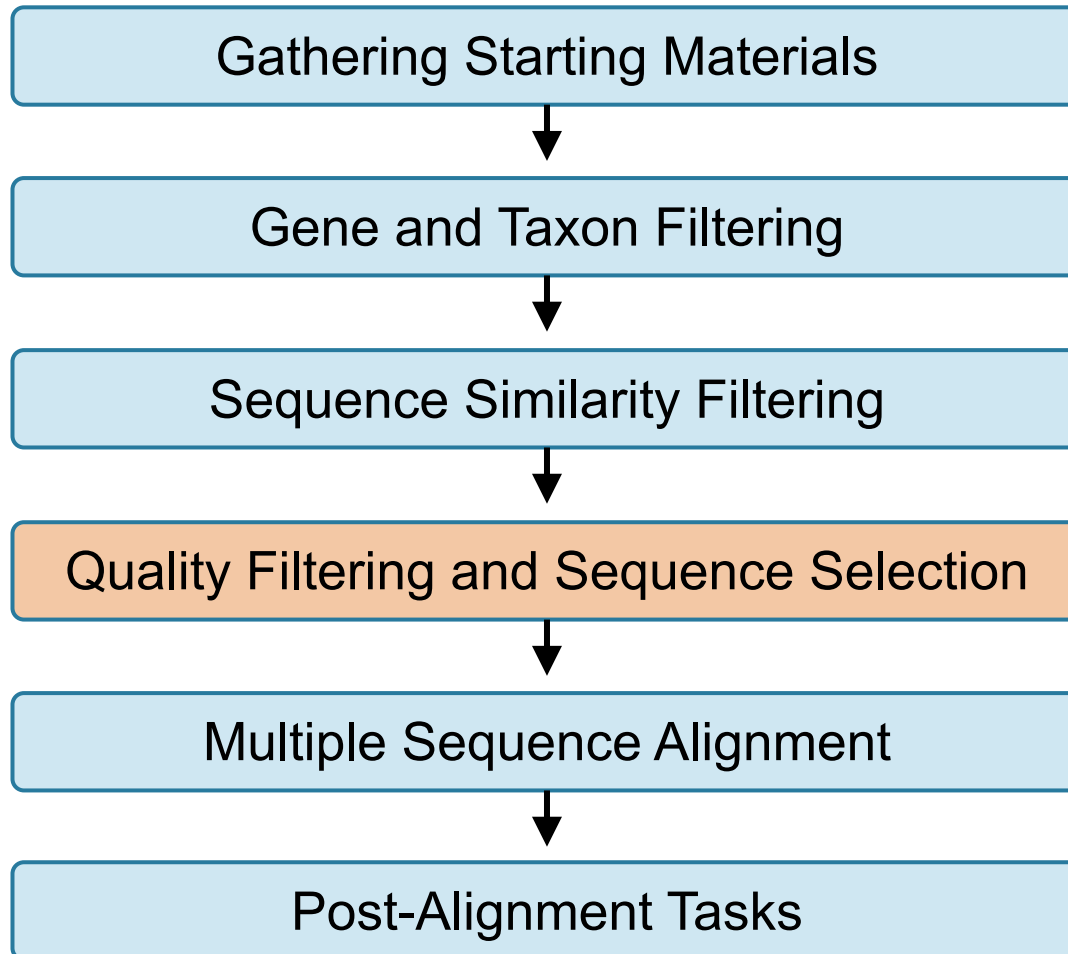
## ‘Simple’ Records

Record	Original Length	Extracted Length	Coordinates Used
GU432436.1	772	772	[0, 771]
GU432437.1	753	753	[0, 752]
GU432438.1	773	773	[0, 772]
GU432439.1	749	749	[0, 748]
JF818222.1	719	717	[2, 718]

## ‘Complex’ Records

Record	Original Length	Extracted Length	Coordinates Used
AM055651.1	510	510	[0, 509]
AM055653.1	508	508	[0, 507]
JF317635.1	17,388	1,575	[1,060, 2,634]
JF317636.1	17,344	1,574	[1,061, 2,634]
AF440085.1	3,700	1,473	[1,014, 2,486]
AF440019.1	3,693	1,470	[1,010, 2,479]
AF338325.1	1,071	41	[1,030, 1,070]
AF338329.1	1,085	84	[1,001, 1,084]

# SuperCRUNCH Workflow





# Quality Filtering and Sequence Selection

---

- What additional filters should be applied to sequences?
- If multiple sequences are available for a species, how do we choose?

# Quality Filtering and Sequence Selection

- What additional filters should be applied to sequences?
- If multiple sequences are available for a species, how do we choose?

## Example from Iguania dataset:

1,426 species

66 loci

58,642 sequences

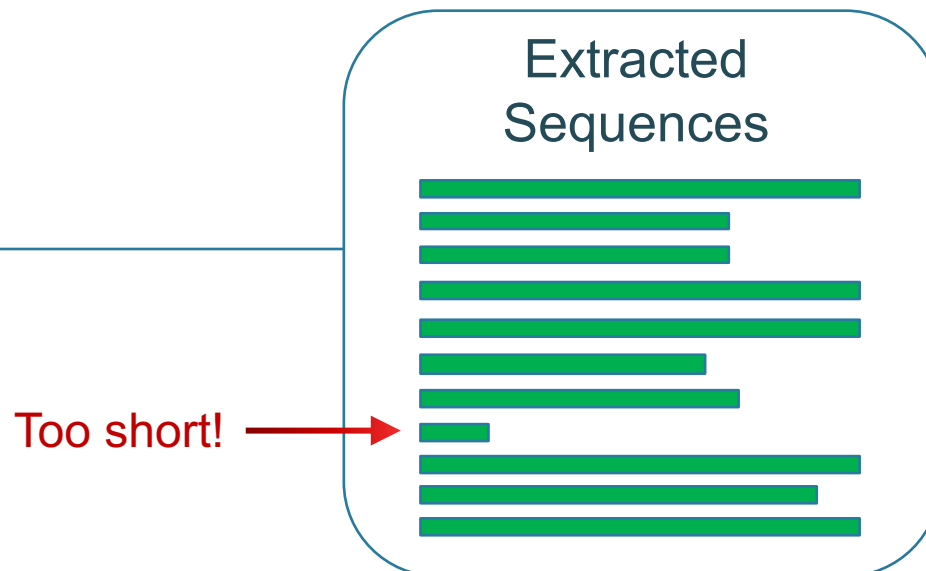


13,419 sequences



# Quality Filtering and Sequence Selection

- What additional filters should be applied to sequences?
  - Minimum length (defined by user)
- If multiple sequences are available for a species, how do we choose?



# Quality Filtering and Sequence Selection

- What additional filters should be applied to sequences?
  - Minimum length (defined by user)
- If multiple sequences are available for a species, how do we choose?
  - 1) Sort by length, select longest sequence
  - 2) Sort by length, test translation (for protein-coding)
  - 3) Random

# Quality Filtering and Sequence Selection

- What additional filters should be applied to sequences?
  - Minimum length (defined by user)
- If multiple sequences are available for a species, how do we choose?
  - 1) Sort by length, select longest sequence
  - 2) Sort by length, test translation (for protein-coding)
  - 3) Random

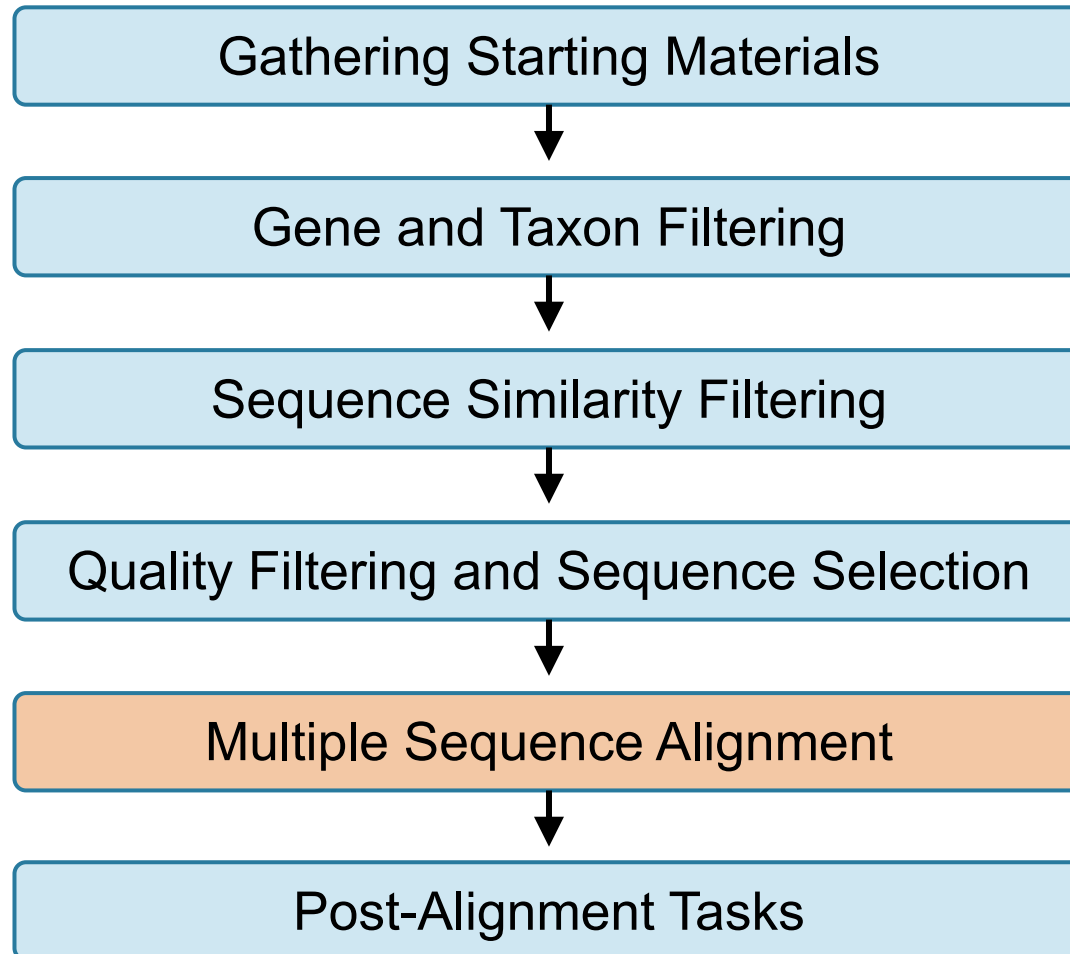
Species-level Data

**One sequence per species per gene**

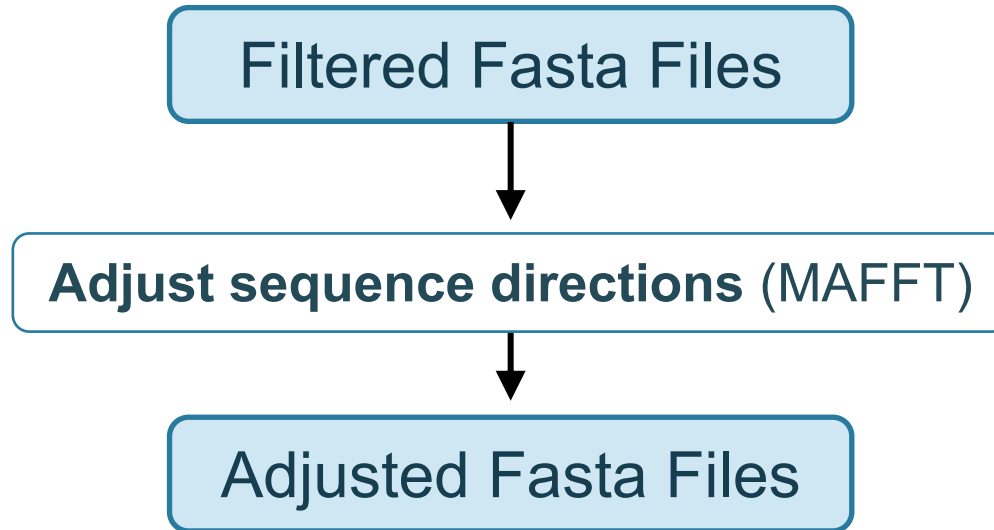
Population-level Data

**All sequences per species per gene**

# SuperCRUNCH Workflow



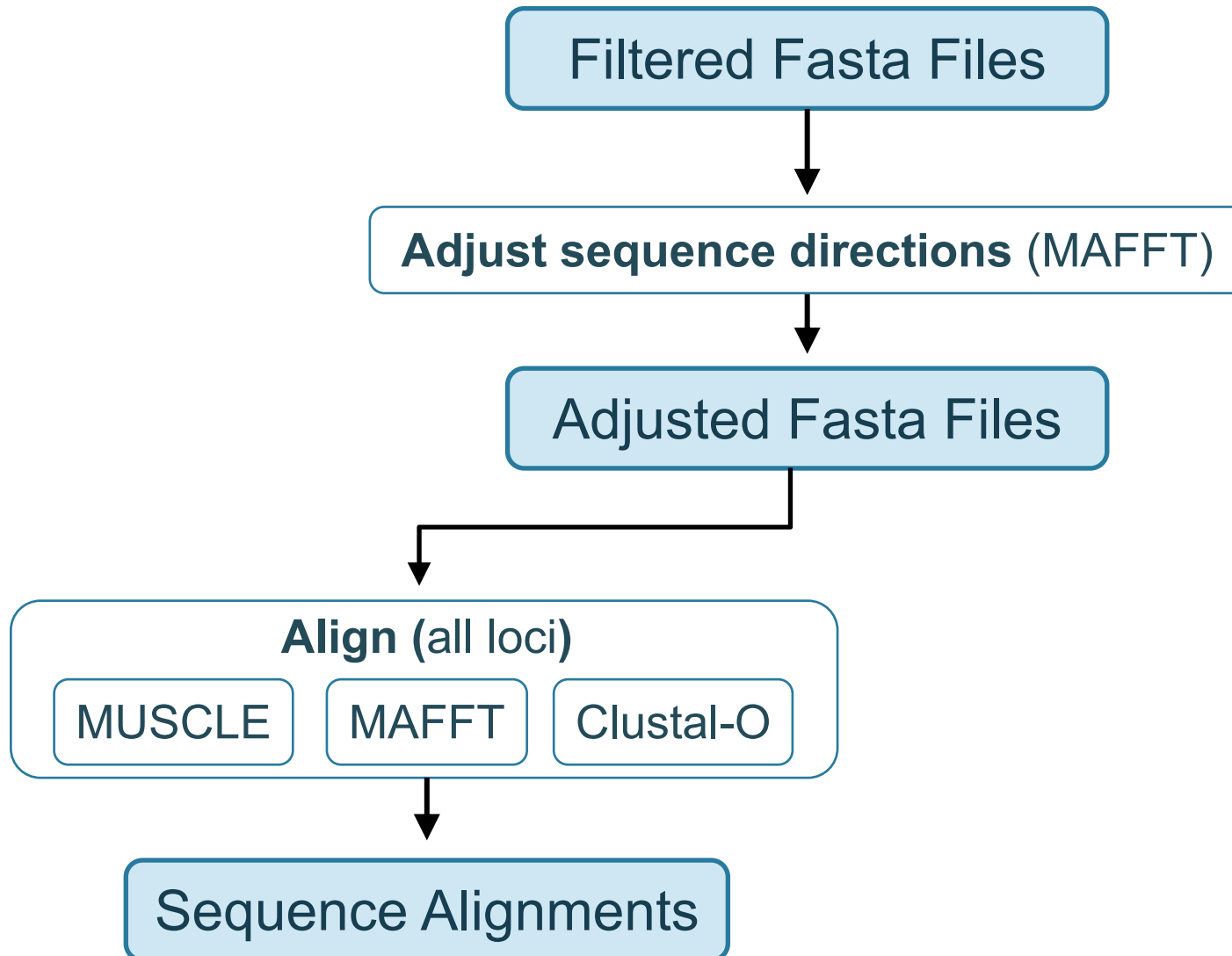
# Multiple Sequence Alignment



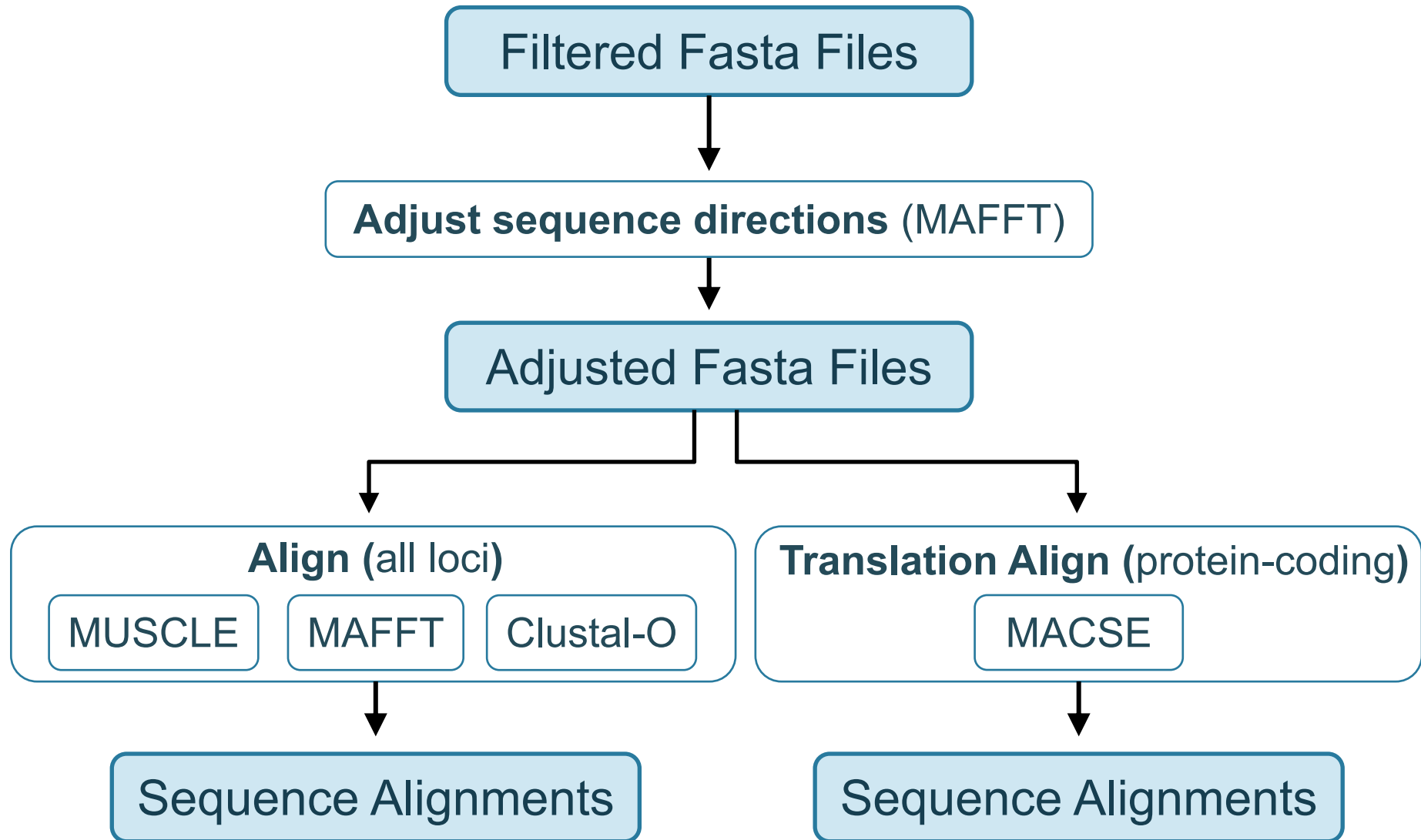
Locus	Correct Direction	Direction Adjusted
UCE-1003	66	5
UCE-1005	62	5
UCE-1012	57	11
UCE-1013	15	4



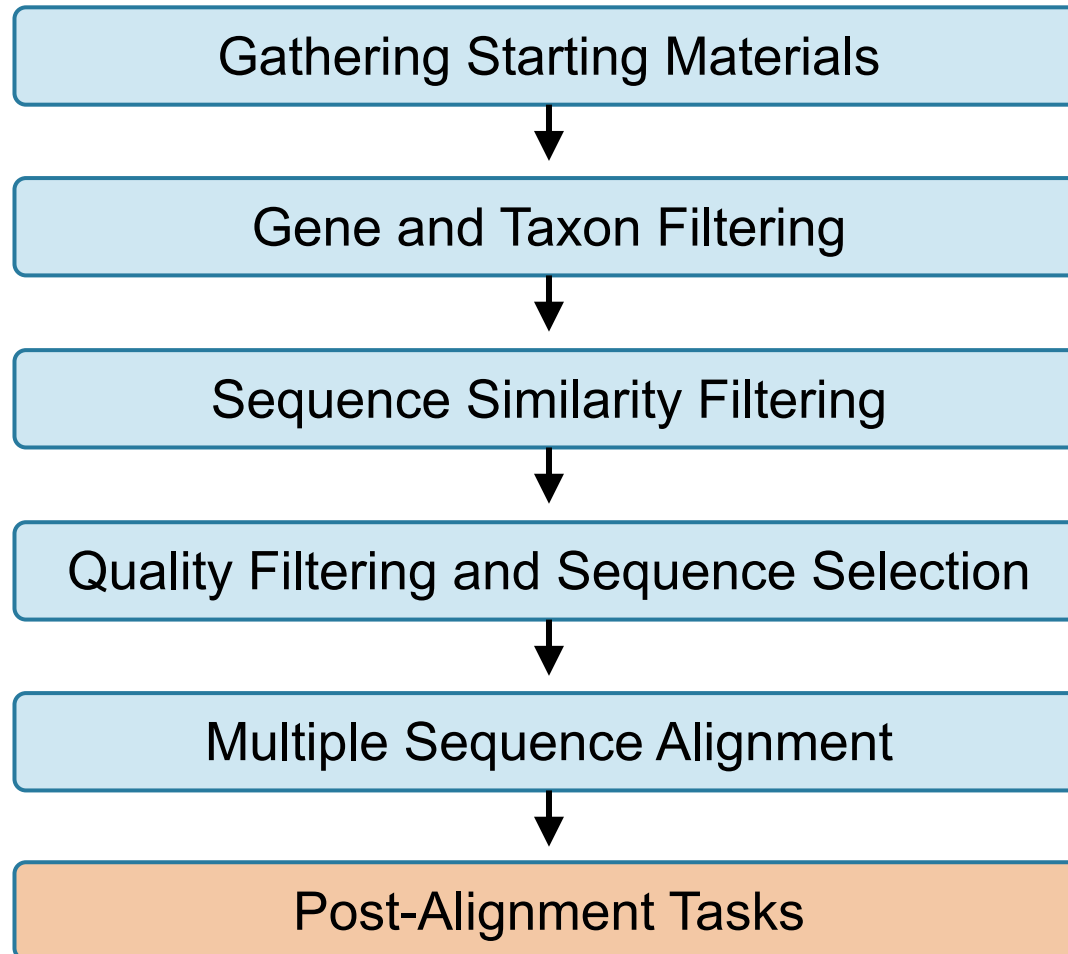
# Multiple Sequence Alignment



# Multiple Sequence Alignment



# SuperCRUNCH Workflow



# Post-Alignment Tasks

---

- Relabel sequence records

# Post-Alignment Tasks

- Relabel sequence records

## Original Labels

```
>JN881132.1 Daboia russelii activity-dependent neuroprotector...  
>KU765220.1 Sceloporus undulatus voucher ADL182 activity-depe...
```

# Post-Alignment Tasks

- Relabel sequence records

## Original Labels

```
>JN881132.1 Daboia russelii activity-dependent neuroprotector...  
>KU765220.1 Sceloporus undulatus voucher ADL182 activity-depe...
```

## Species

```
>Daboia_russelii  
>Sceloporus_undulatus
```

# Post-Alignment Tasks

- Relabel sequence records

## Original Labels

```
>JN881132.1 Daboia russelii activity-dependent neuroprotector...  
>KU765220.1 Sceloporus undulatus voucher ADL182 activity-depe...
```

## Species

```
>Daboia_russelii  
>Sceloporus_undulatus
```

## Accession

```
>JN881132.1  
>KU765220.1
```

# Post-Alignment Tasks

- Relabel sequence records

## Original Labels

```
>JN881132.1 Daboia russelii activity-dependent neuroprotector...  
>KU765220.1 Sceloporus undulatus voucher ADL182 activity-depe...
```

## Species

```
>Daboia_russelii  
>Sceloporus_undulatus
```

## Accession

```
>JN881132.1  
>KU765220.1
```

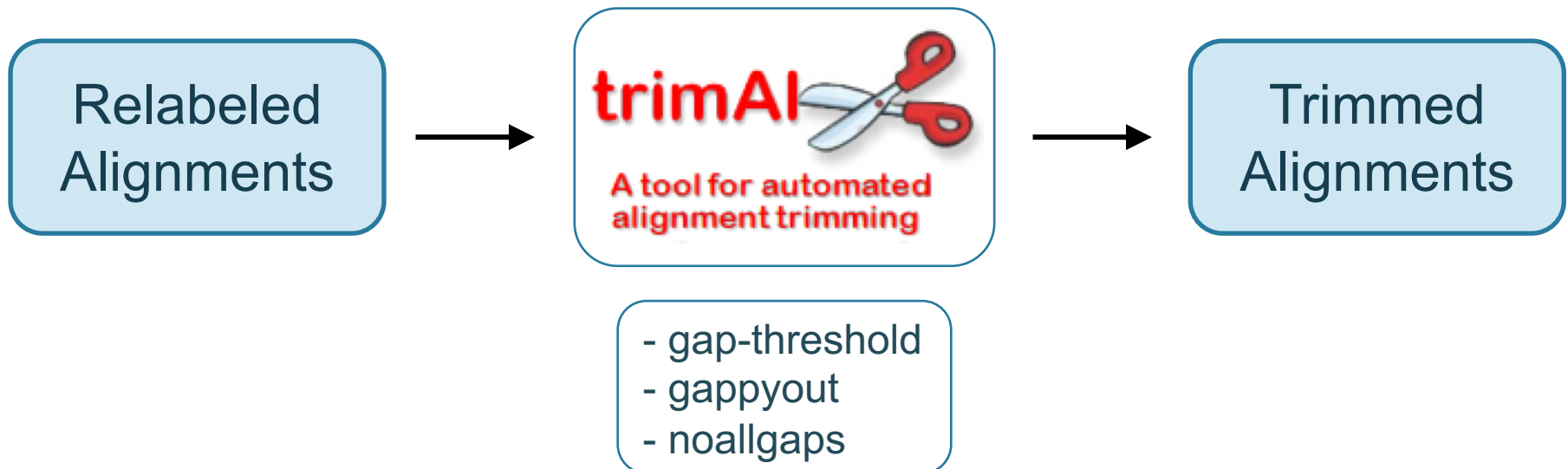
## Species + Accession

```
>Daboia_russelii_JN881132.1  
>Sceloporus_undulatus_KU765220.1
```



# Post-Alignment Tasks

- Relabel sequence records
- Trim alignments



# Post-Alignment Tasks

- Relabel sequence records
- Trim alignments
- Convert to nexus, phylip format

# Post-Alignment Tasks

- Relabel sequence records
- Trim alignments
- Convert to nexus, phylip format
- Concatenate alignments

## Requires Species Labeling

```
>Daboia_russelii  
>Sceloporus_undulatus
```

## Missing Data Options

N - ?

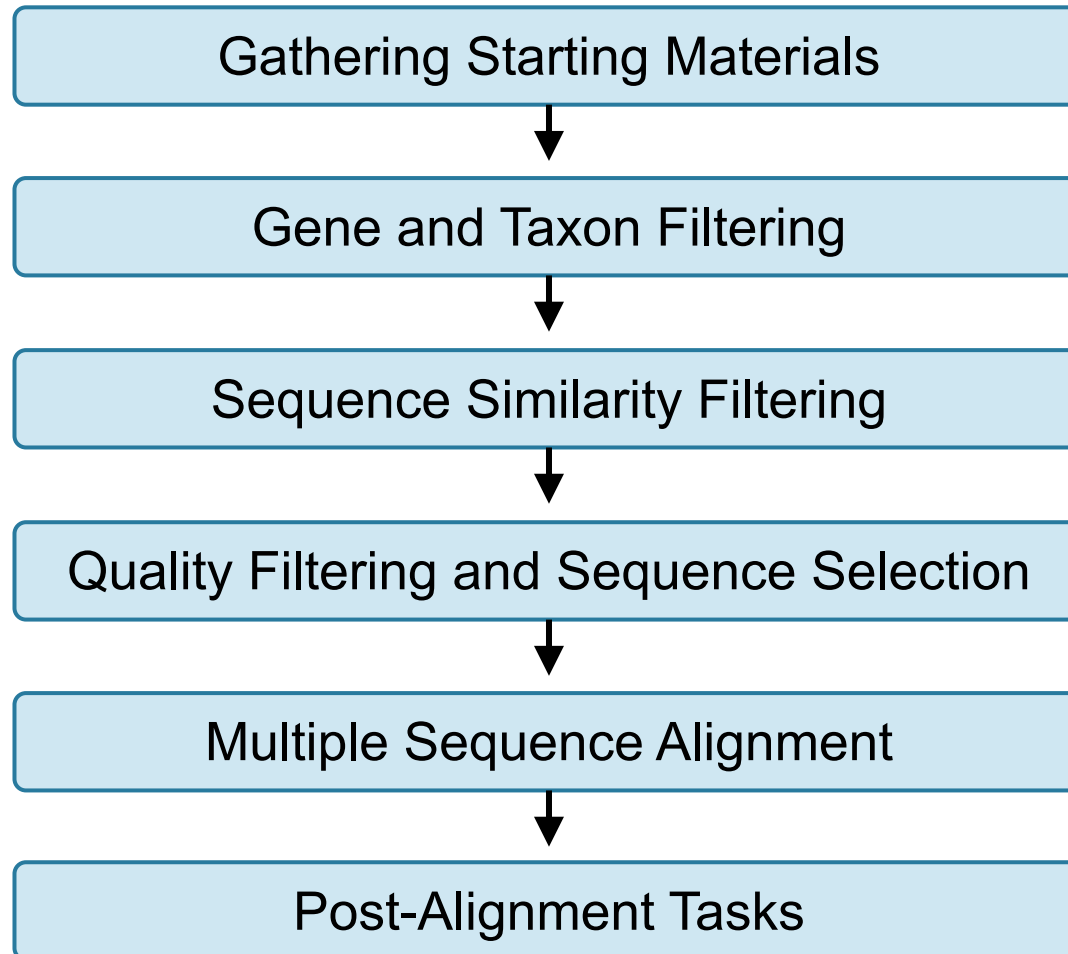
## Outputs

Fasta or phylip

Partitions file

Loci per taxon

# SuperCRUNCH Workflow



# SuperCRUNCH Workflow

Gathering Starting Materials



Post-Alignment Tasks

# Benchmarking: Time

- Dataset dependent, but reasonable!
- Bottleneck is always sequence alignment

	Population level	Traditional Supermatrix	Genomic (UCE's)
File Size	52 MB	13.3 GB	205 MB
Sequences	82,557	8,785,378	338,942
Taxa	3	1,426	123
Genes	4	70	5,041
Total time	~2.5 min	~4.5 hours	~10.5 hours

# Benchmarking: Other Programs

## Comparison to PyPHLAWD

- PyPHLAWD uses GenBank release database and Taxonomy Browser to fetch starting sequences
- Default is automated clustering of all sequences (requires follow-up identification and processing)
- A 'baited analysis' can target specific loci, requires reference sequences

# Benchmarking: Other Programs

## Comparison to PyPHLAWD

- Compared supermatrix construction for Iguania, large clade of squamate reptiles (~1900 species)
- Targeted 68 loci
  - 61 nuclear – used automatic reference selection
  - 7 mtDNA loci – used validated reference sequences

### Three analyses:

**SuperCRUNCH**  
NCBI direct-download

**SuperCRUNCH**  
PyPHLAWD database

**PyPHLAWD**  
PyPHLAWD database



# Benchmarking: Other Programs

## Comparison to PyPHLAWD

	Loci	Taxa	Sequences
<b>SuperCRUNCH</b> NCBI direct-download	66	1,426	13,419
<b>SuperCRUNCH</b> PyPHLAWD database	66	1,358	12,654
<b>PyPHLAWD</b> PyPHLAWD database	65	1,069	10,397

# Benchmarking: Other Programs

- PyPHLAWD was comparable for nuclear loci
- PyPHLAWD failed to find 62% (2,125) of mtDNA seqs

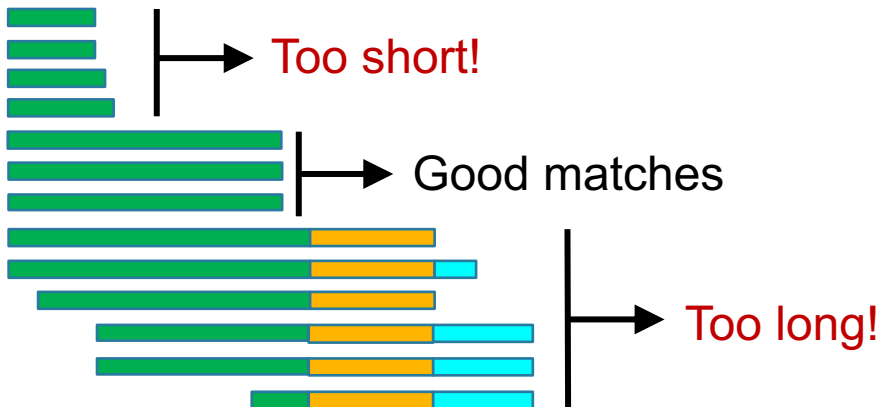
# Benchmarking: Other Programs

- PyPHLAWD was comparable for nuclear loci
- PyPHLAWD failed to find 62% (2,125) of mtDNA seqs

## References



Sequences must be highly similar to bait to pass filter  
Entire sequence is kept or discarded (no extraction)



# Looking for Testers!

## SuperCRUNCH FOR PHYLOGENETIC DATA



Pre-Print

<https://doi.org/10.1101/538728>

Download

<https://github.com/dportik/SuperCRUNCH>

Documentation

<https://github.com/dportik/SuperCRUNCH/wiki>

Tutorials & Data

<https://osf.io/bpt94/>



# Benchmarking

	Population level	Genomic (UCE's)	Supermatrix
File Size	52 MB	205 MB	13.3 GB
Sequences	82,557	338,942	8,785,378
Taxa	3	123	1,426
Genes	4	5,041	70
Gene and Taxon Filtering	5s	2hr 40min	34min
Sequence Similarity Filtering	8s	---	1hr 10min
Sequence Selection	1s	7min	7min
Sequence Alignment	2min	7hr 10min	2hr 40min
Post-Alignment Tasks	2s	21min	7s
Total cpu time	<b>~2.5 min</b>	<b>~10.5 hours</b>	<b>~4.5 hours</b>