# New Sequencing Technologies and Hybrid Assemblies

Harindra M. Arachchi, Manuel Garber, Chad Nusbaum
Sarah Young, Michael C. Zody, David Jaffe, Michael Fitzgerald

# Objective

- Obtain PCR template from unclonable gaps
- Take advantage of new technologies to sequence them
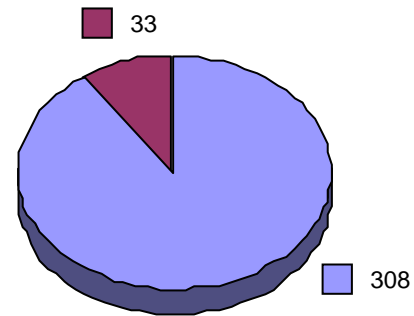- Develop methodology for closing these gaps

# Finished human genome (2004)

Finished human genome build 35

341 gaps

33

308

Heterochromatic

Euchromatic

2.85 billion nucleotides

(99% of the euchromatic genome)

Unclonable

Segmental duplications

# Gaps on human chromosome 15

10 gaps

→ 7 in segmental duplicated regions

→ 3 apparently unclonable

→ Not captured in human tiling path

→ Not captured in chimpanzee tiling path

53x physical coverage screened

# Overview.

- Process
- Analysis
- Conclusion

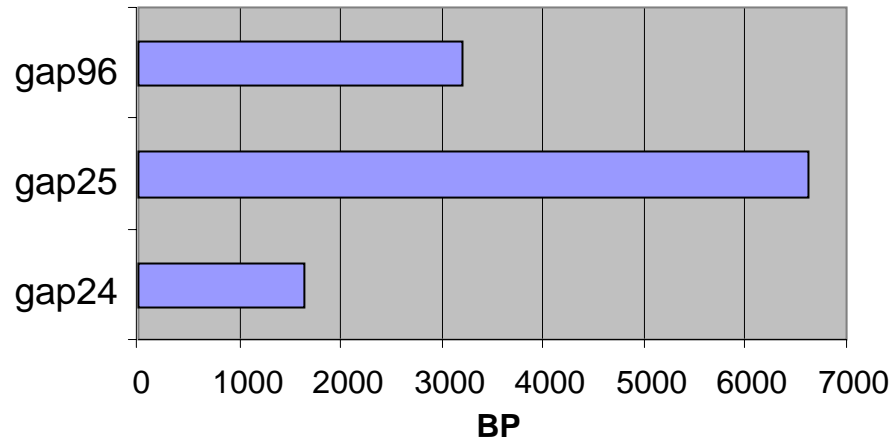# Overview.

- **Process**
- Analysis
- Conclusion

# STEP 1: Align HGP & Celera

gap

HGP

CELERA

Celera coverage of gap

gap96

gap25

gap24

0    1000   2000   3000   4000   5000   6000   7000

**BP**

# STEP 2: Integrate shatter reads

gap

HGP

CELERA

SHATTER

Shatter library coverage of gap

gap96

gap25

gap24

| | 0 | 1000 | 2000 | 3000 | 4000 | 5000 | 6000 | 7000 |

**BP**

# STEP 3: Integrate 454 data

454 sequence

Assemble 454

Extract contig consensus

Map quality

Create artificial reads(AR)

Integrate with Sanger

Integrated

gap

HGP

CELERA

SHATTER

454 AR

454 coverage of gap

gap96

gap25

gap24

0          5000        10000       15000

BP

□ Automated assembly
■ Manual extension

# Newbler assembly breaks at SNPs

```
CTCTGCTCATTTCAGCTCGGACGGTGGTCCCTT
CTCTGCTCATTTCAGCTCGGACGGTGGTCCCTTAAGCAGGCCGAAACTGATGGTCTCATCTCCTGCACGCTC        Sanger data
CTCTGCTCATTTCAGCTCGGACGGTGGTCCCTTCAGCAGGCCGAAACTGATGGTCTCATCTCCTGCACGCTC
                                 AGCAGGCCGAAACTGATGGTCTCATCTCCTGCACGCTC        454 data
```

Haplotype difference

# 454 Assembly: ALLPATHS

ALLPATHS assembly algorithm

-Finds all shared kmers between reads

-Uses shared kmers to build graph across the data

-Graphs represents contigs as edges, branch points as nodes



David Jaffe

# 454 Assembly: ALLPATHS

ALLPATHS assembly algorithm

-Finds all shared kmers between reads

-Uses shared kmers to build graph across the data

-Graphs represents contigs as edges, branch points as nodes



Complete path across region: spans a gap

David Jaffe

# 454 Assembly: ALLPATHS

ALLPATHS assembly algorithm

-Finds all shared kmers between reads

-Uses shared kmers to build graph across the data

-Graphs represents contigs as edges, branch points as nodes



Branched path = two haplotypes

David Jaffe

# 454 Assembly: ALLPATHS

ALLPATHS assembly algorithm

-Finds all shared kmers between reads

-Uses shared kmers to build graph across the data

-Graphs represents contigs as edges, branch points as nodes

PCR slippage in repeat

David Jaffe
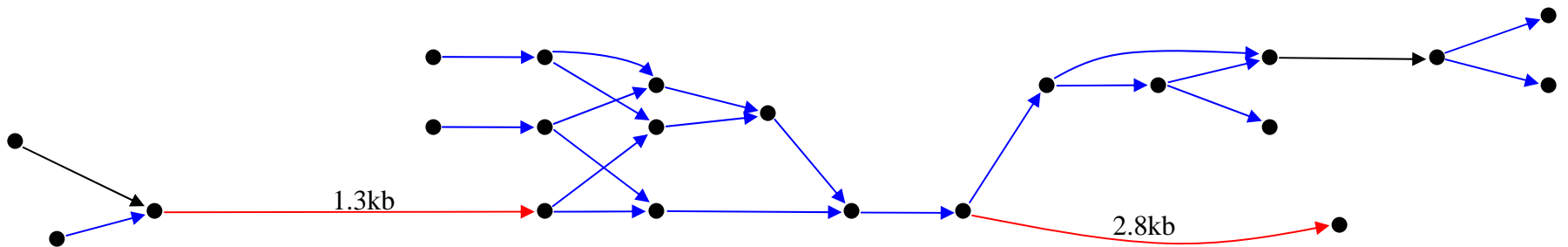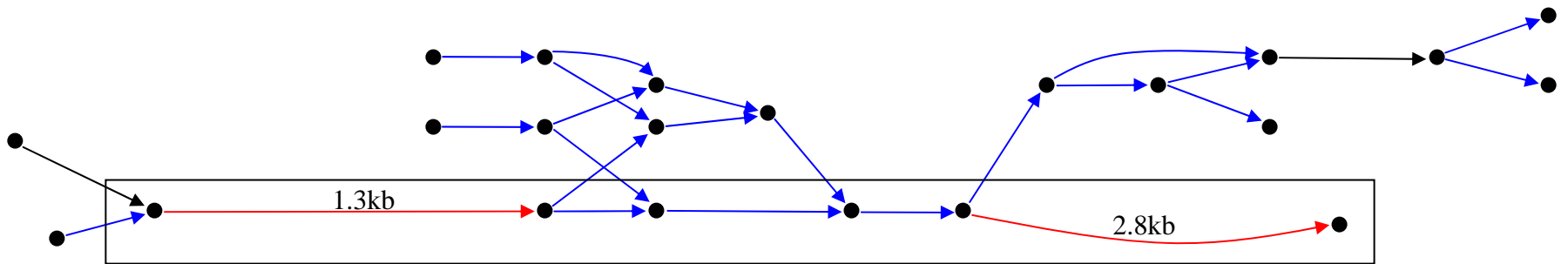
# 454 Assembly: ALLPATHS

ALLPATHS assembly algorithm

-Finds all shared kmers between reads

-Uses shared kmers to build graph across the data

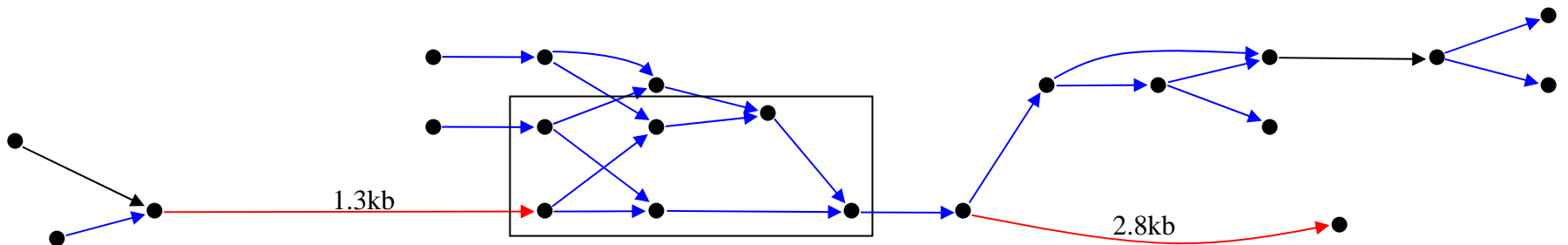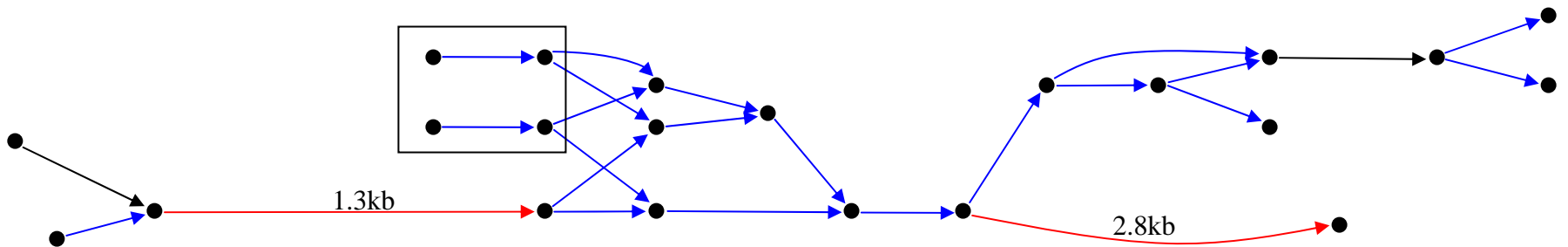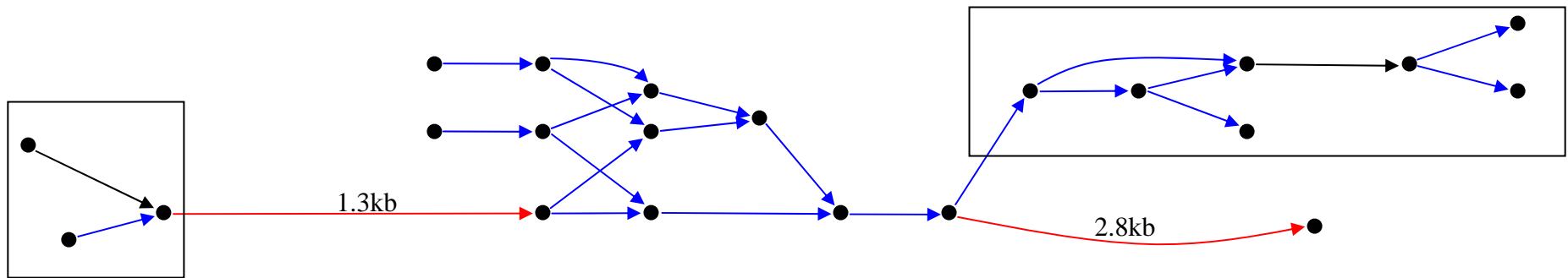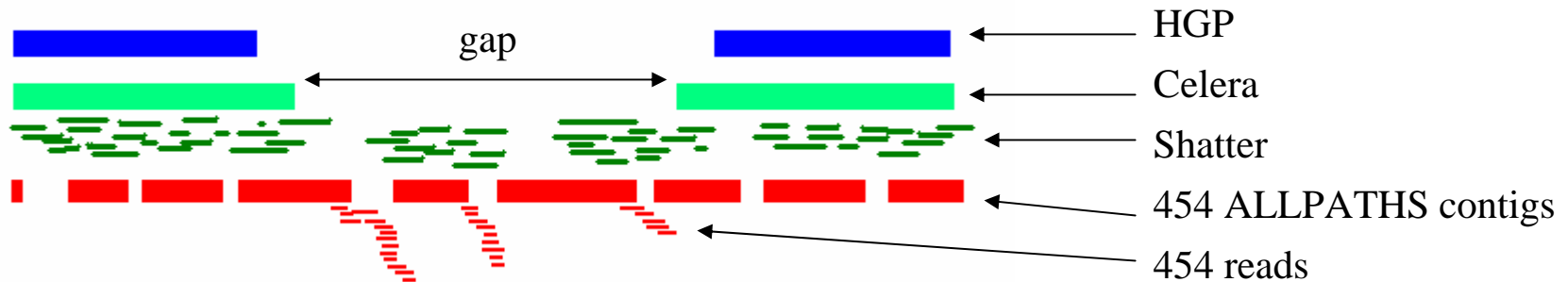-Graphs represents contigs as edges, branch points as nodes



Artifacts from messy PCR

David Jaffe

# Manual extension of contigs



gap
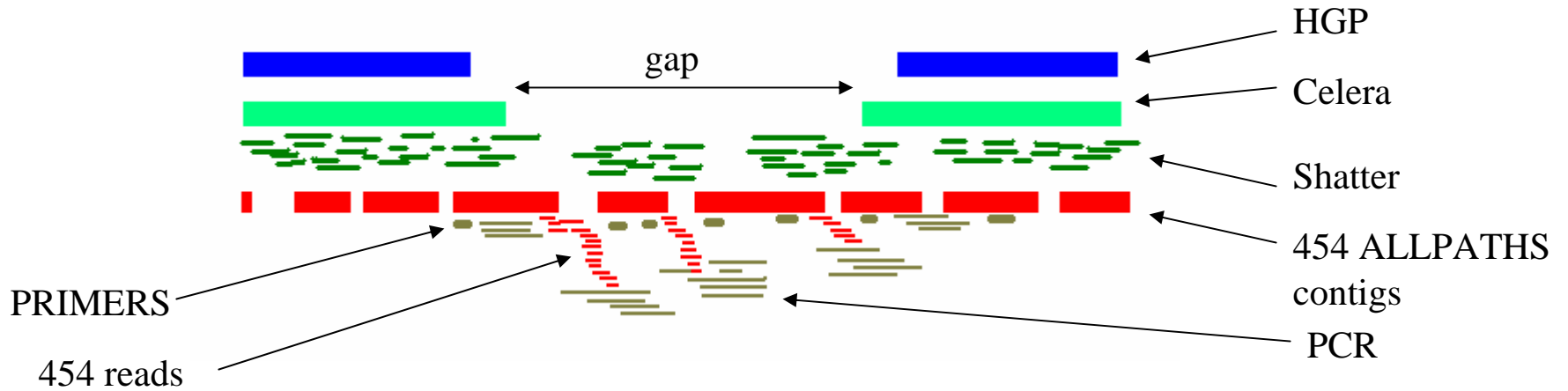
HGP

Celera

Shatter

454 ALLPATHS contigs

454 reads

- Identify unassembled reads computationally by string search

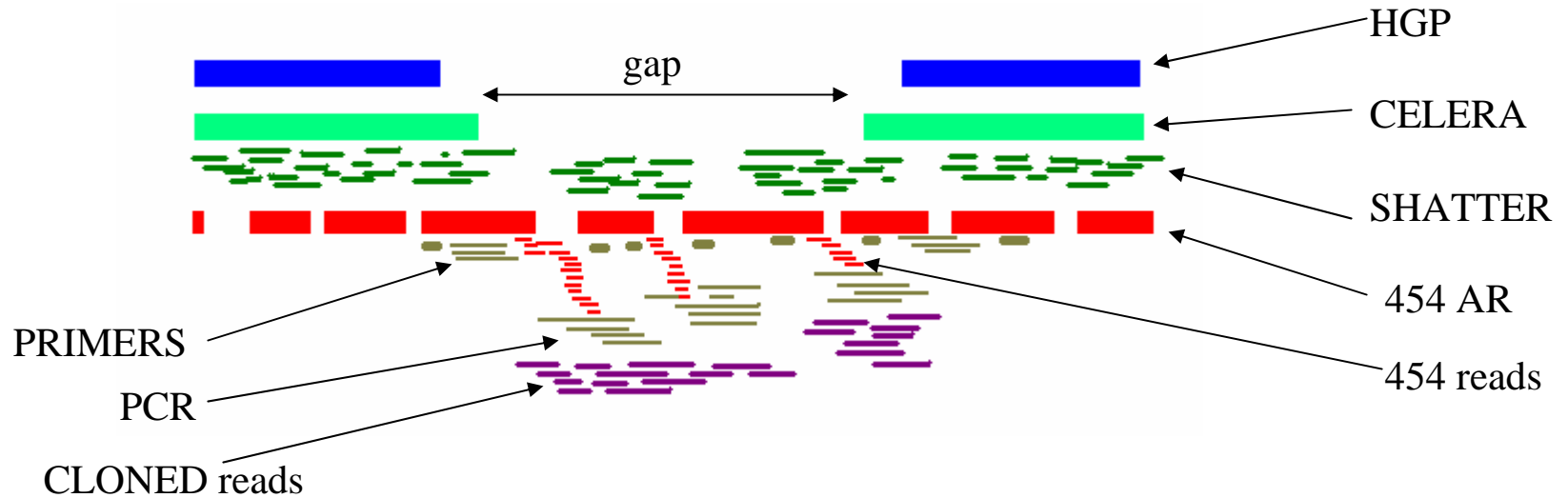- Manually insert reads into assembly

454 contig    GTCT

"Fished" in 454 reads    ATGT

```
GTACGTGTTACATGTCTGATGTATGAGGTGTGTGTGGTACGTGTGTTATATGC
 GTTACATGTCTGATGTATGAGGTGTGTGTGGTACGTGTGTTATATGCAACATGTGTGGTGTAGGCAATGTGT
  ACATGTCTGATGTATGAGGTGTGTGTGGTACGTGTGTTATATGCAACATGTGTGGTGTAGGCAATGTGT
   ATGTCTGATGTATGAGGTGTGTGTGGTACGTGTGTTATATGCAACATGTGTGGTGTAGGCAATGTGT
   ATGTCTGATGTATGAGGTGTGTGTGGTACGTGTGTTATATGCAACATGTGTGGTGTAGGCAATGTGT
    TGTCTGATGTATGAGGTGTGTGTGGTACGTGTGTTATATGCAACATGTGTGGTGTAGGCAATGTGT
     GTCTGATGTATGAGGTGTGTGTGGTACGTGTGTTATATGCAACATGTGTGGTGTAGGCAATGTGT
      TCTGATGTATGAGGTGTGTGTGGTACGTGTGTTATATGCAACATGTGTGGTGTAGGCAATGTGT
       CTGATGTATGAGGTGTGTGTGGTACGTGTGTTATATGCAACATGTGTGGTGTAGGCAATGTGT
           AGGTGTGTGTGGTACGTGTGTTATATGCAACATGTGTGGTGTAGGCAATGTGT
           AGGTGTGTGTGGTACGTGTGTTATATGCAACATGTGTGGTGTAGGCAATGTGT
           AGGTGTGTGTGGTACGTGTGTTATATGCAACATGTGTGGTGTAGGCAATGTGT
             TGTGTGGTACGTGTGTTATATGCAACATGTGTGGTGTAGGCAATGTGT
                TGTTATATGCAACATGTGTGGTGTAGGCAATGTGT
                TTATATGCAACATGTGTGGTGTAGGCAATGTGT
                 TATATGCAACATGTGTGGTGTAGGCAATGTGT
                  ATATGCAACATGTGTGGTGTAGGCAATGTGT
                   TTATGCAACATGTGTGGTGTAGGCAATGTGT
                     CAACATGTGTGGTGTAGGCAATGTGT
                     CAACATGTGTGGTGTAGGCAATGTGT
                      CAACATGTGTGGTGTAGGCAATGTGT
                       AACATGTGTGGTGTAGGCAATGTGT
                       AACATGTGTGGTGTAGGCAATGTGT
                          TGTGGTGTAGGCAATGTGT
                           TGGTGTGGGCAATGTGT
                            GGTGTAGGCAATGTGT
                             GTAGGCAATGTGT
                             GTAGGCAATGTGT
```

# STEP 4: PCR to verify 454

# STEP 5: Cloning PCR fragments

# Overview

- Process
- **Analysis**
- Conclusion

# Hard finishing….

- Sequence flanking all 3 gaps looks normal
- Problem is nasty repeat in gap
  - Unable to clone
  - Can sequence by 454 but Newbler can't assemble
  - ALLPATHS gave better assembly
  - Had to finish the job by hand

# Gap25

Status: Closed

Size: 10225bp

Unclonable: 21.6%

Total 454 coverage: 98%

Repeat motif:

GGTGTTTGTGTGTATGGT

# Gap96

Status: Closed

Size: 5474bp

Unclonable: 14%

Total 454 coverage: 83.7%

Repeat motif: TATGTGTGTGGCATGTGGT

# Gap24

Status: Active

Deletions complications

Size: 2539bp**

Unclonable: 14%

Total 454 coverage: 92%

Repeat motif:

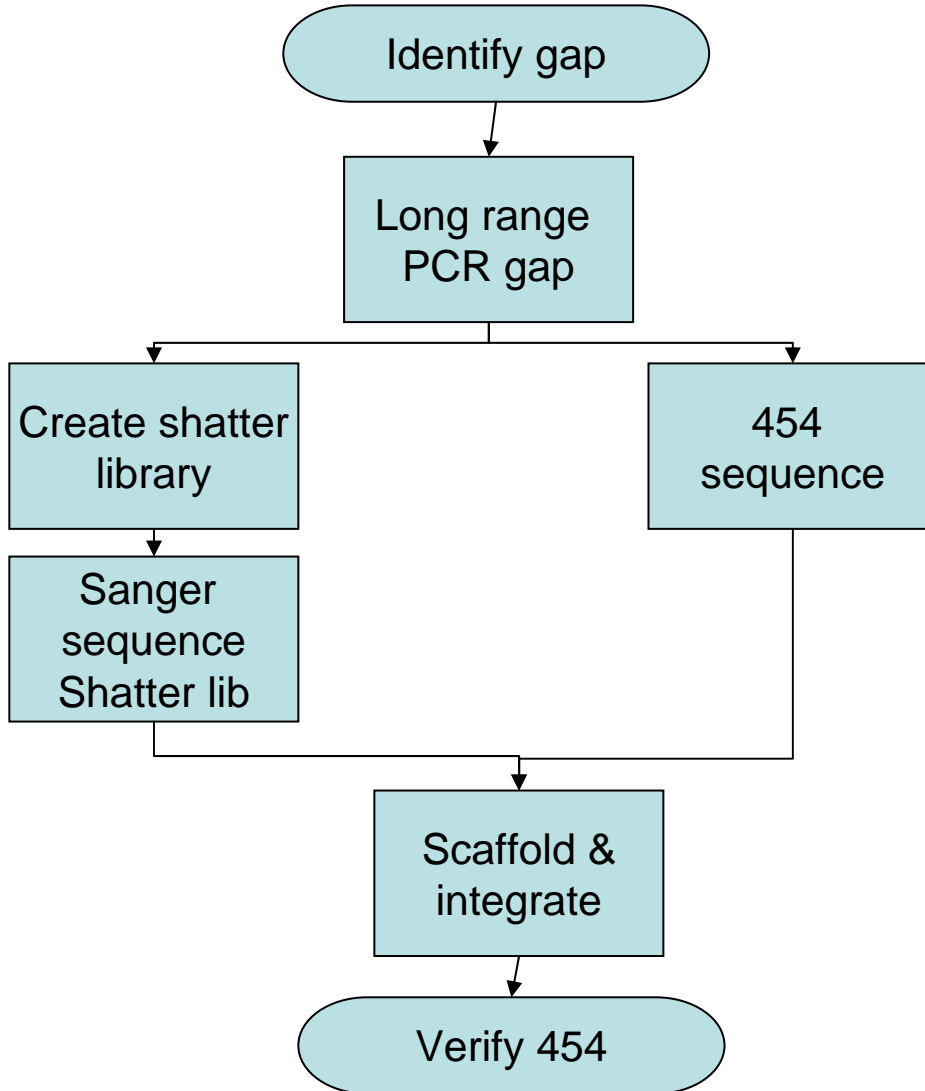TGTATGGTGTGTGGCGTGTG

**sequence captured so far

# Overview

- Process
- Analysis
- **Conclusion**

# Summary

HGP unclonable remaining:126

- Not due to copy number variations
- In unique regions

```
        ┌─────────────────┐
        │   Identify gap   │
        └─────────────────┘
                 │
                 ▼
        ┌─────────────────┐
        │   Long range     │
        │   PCR gap        │
        └─────────────────┘
          │             │
          ▼             ▼
  ┌──────────────┐  ┌──────────────┐
  │ Create shatter│  │     454       │
  │   library     │  │  sequence     │
  └──────────────┘  └──────────────┘
         │                  │
         ▼                  │
  ┌──────────────┐          │
  │   Sanger      │          │
  │  sequence     │          │
  │ Shatter lib   │          │
  └──────────────┘          │
         │                  │
         └────────┬─────────┘
                  ▼
        ┌─────────────────┐
        │   Scaffold &     │
        │   integrate      │
        └─────────────────┘
                 │
                 ▼
        ┌─────────────────┐
        │   Verify 454     │
        └─────────────────┘
```

How do we define "finished" ?

# Thank you!

Manuel Garber
Sarah Young
Chad Nusbaum
Michael C. Zody

Broad Institute
Special Projects Group
Finishing Group
Genome Sequencing Platform

Michael Fitzgerald
Will Lee
David Jaffe
Lisa Zembek
Niall Lennon

Cristyn Kells
James Bonfield