# Sequencing the *Phaseolus vulgaris* G19833 Genome

Scott Jackson, Purdue University

Phil McClean, North Dakota State University

Perry Cregan, USDA-ARS, Beltsville

Dan Rokhsar, DOE Joint Genome Institute

Jeremy Schmutz, HudsonAlpha Institute

David Hyten, USDA-ARS, Beltsville

# Project overview

- Produce as good as possible reference sequence for PV with available funds that will serve as a directly useful tool for common bean improvement (when combined with ongoing Bean-CAP efforts) and a reference comparator for soybean

- The DOE Joint Genome Institute is producing a 454/Illumina based draft sequence for soybean comparative work

- USDA has funded the project "**A sequence map of the common bean genome for bean improvement"** to add longer range linking data and genetic mapping data to the draft in order to produce a reference PV genome sequence

- The overlap between the DOE-JGI portion and the USDA portion is coordinated by Jeremy Schmutz at HudsonAlpha, who is also funded as part of the JGI Plant Genome Program

# Expected genome data for PV V1.0

- 12x 454 linear data

- 4-6x 454 paired data from multiple sized libraries (4kb and 8kb)

- 50x of Illumina data

- 200k fosmid Sanger sequenced pairs

- 80k BAC Sanger sequenced pairs (includes 40k from BeanMap efforts)

- Up to 1400 marker genetic map (Hyten & Cregan)

- BAC physical map (phaseolus.genomics.purdue.edu)

# 454 linear sequencing

- Collected 7.1 GB of linear data (~11x HQ bases)

- Includes 26 runs from 4 libraries

- Average HQ read length of 274 bps

```
Phred20     Reads  %ofReads
  0- 49    648713   2.5%    |X
 50- 99   2134679   8.2%    |XXXX
100-149   2073312   7.9%    |XXXX
150-199   2560947   9.8%    |XXXXX
200-249   3063228  11.7%    |XXXXXX
250-299   3561945  13.7%    |XXXXXXX
300-349   4114373  15.8%    |XXXXXXXX
350-399   3769733  14.4%    |XXXXXX
400-449   2479550   9.5%    |XXXXX
450-499   1418195   5.4%    |XXX
500-549    264387   1.0%    |X
550-599      4896   0.0%    |
600-649        19   0.0%    |
650-699         8   0.0%    |
700-749         3   0.0%    |

Number of reads: 26,093,988
Total bases: 8,881,905,934
Total Phred 20 bases: 7,159,549,663
Average length: 340.4
Phred average: 274.4
Phred average without failures: 276.9
Percent failed: 1.0%
```

# 454 paired sequencing

- Collected 3.4 Gb of total HQ raw data

- Includes 14 runs from 8 libraries

- Once true pairs are determined we only collect 724 Mb or ~1.1x of paired data

- Making good paired libraries has been the most difficult part of this project

- We have 4 more library attempts in progress

```
Phred20     Reads %ofReads
  0- 49   338969   6.6%   |XXX
 50- 99  1506543  29.4%   |XXXXXXXXXXXXXX
100-149  1218194  23.8%   |XXXXXXXXXXX
150-199   923676  18.0%   |XXXXXXXXX
200-249   601244  11.7%   |XXXXX
250-299   329828   6.4%   |XXX
300-349   145060   2.8%   |X
350-399    47794   0.9%   |
400-449     6597   0.1%   |
450-499       87   0.0%   |

Number of reads: 5,117,992
Total bases: 840,722,607
Total Phred 20 bases: 723,932,906
Average length: 164.3
Phred average: 141.4
Phred average without failures: 145.2
Percent failed: 3.3%
```

# Illumina sequencing

- We have collected several different GA2x flow cells for PV, most of the reads are paired 76mers that will be used to correct 454 read errors

```
Phred20        Reads %ofReads
 0- 9   15837191   2.9%     |X
10-19    6520659   1.2%     |X
20-29    8564932   1.6%     |X
30-39   12494229   2.3%     |X
40-49   18842458   3.4%     |XX
50-59   31401672   5.7%     |XXX
60-69   71041097  12.9%     |XXXXXX
70-79 383902008  70.0%     |XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX

Number of reads: 548,604,246
Total bases: 41,693,922,696
Total Phred 20 bases: 36,746,552,154
Average length: 76.0
Phred average: 67.0
Phred average without failures: 71.2
Percent failed: 7.9%
```

# Fosmid libraries

- Goal is to produce about 400k reads or 200k pairs from fosmid length libraries

- 2 libraries have been constructed
  - VUL    33,465 bp insert; 200 plates
  - VUK    35,558 bp insert; 260 plates

- We have sequenced a combined 152k reads, 100k more are in progress now at HudsonAlpha

- We are sequencing more BES from the new BAC library to replace some of the FES to beef up long range linking

# New BAC library

- Dave Kudrna at Arizona Genomics created a new BAC library called PV_GBb with the goal of maximizing the insert length

- The first section is 120 plates with 149kb average

- The second section is 105 plates with a 136kb insert

- We are sequencing the entire first section and filling in additional inserts from the second section

- These will add to the 89K reads from PV_Gba and give us decent BAC coverage across the genome

| Ligations picked | into Lib. Plates # | QC plate Information | avg inserts |
|---|---|---|---|
| A2-1 | Plates 1 - 41 | QC Plates 1 & 2: (Lib plts 1- 95) | 148.4 kb (plates 1 - 120) |
| A2-2-1 | Plates 42 - 46 | | |
| A2-1-1 | Plates 47 - 95 | | |
| A2-2-1 | Plates 96 - 120 | QC Plate 3: (Lib plt 96 - 120) | |

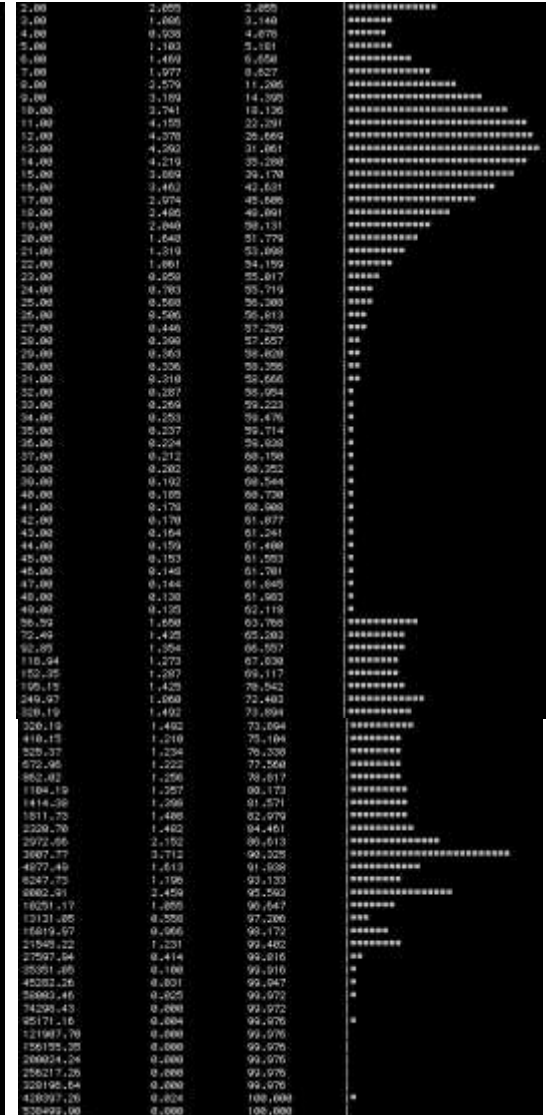| A2-2-1 | Plates 121 - 141 | QC Plate 4: (lib plt 121 - 216) | 136 kb (plates 121 - 216) |
|---|---|---|---|
| B-B | Plates 142 - 226 | | |

# GC profiles of data sets

# Kmer distributions



Sanger

454 Pairs

454 Linear

# Initial Newbler assemblies 454 only

- ~8x assembly

```
Main genome scaffold total: 20067
Main genome contig total:    61136
Main genome scaffold sequence total: 419.1 MB
Main genome contig sequence total:    367.2 MB (—> 12.4% gap)
Main genome scaffold N/L50: 2046/55.9 KB
Main genome contig N/L50:    11467/9.1 KB
Number of scaffolds > 50 KB: 2385
% main genome in scaffolds > 50 KB: 54.3%
```

- ~12x assembly

```
Main genome scaffold total: 10037
Main genome contig total:    47131
Main genome scaffold sequence total: 470.1 MB
Main genome contig sequence total:    416.5 MB (—> 11.4% gap)
Main genome scaffold N/L50: 520/214.4 KB
Main genome contig N/L50:    7982/14.8 KB
Number of scaffolds > 50 KB: 1998
% main genome in scaffolds > 50 KB: 82.4%
```

- These builds are the basis for the genetic mapping effort

| Minimum Scaffold Length | Number of Scaffolds | Number of Contigs | Total Scaffold Length | Total Contig Length | Scaffold Contig Coverage |
|---|---|---|---|---|---|
| All | 10,037 | 47,131 | 470,076,288 | 416,492,511 | 88.60% |
| 1 kb | 10,037 | 47,131 | 470,076,288 | 416,492,511 | 88.60% |
| 2.5 kb | 8,010 | 45,104 | 465,583,552 | 411,999,785 | 88.49% |
| 5 kb | 5,753 | 42,714 | 458,127,389 | 404,765,693 | 88.35% |
| 10 kb | 4,525 | 40,552 | 449,732,354 | 397,781,006 | 88.45% |
| 25 kb | 3,071 | 36,351 | 425,878,517 | 378,081,319 | 88.78% |
| 50 kb | 1,998 | 31,259 | 387,529,708 | 346,509,091 | 89.41% |
| 100 kb | 1,157 | 24,832 | 328,052,126 | 295,902,874 | 90.20% |
| 250 kb | 427 | 14,827 | 213,547,125 | 195,179,721 | 91.40% |
| 500 kb | 145 | 7,558 | 115,721,896 | 106,708,733 | 92.21% |
| 1 mb | 22 | 2,025 | 33,505,576 | 31,309,850 | 93.45% |

# Latest Newbler assembly

- We are running several version of the compete data set to date, including BES and FES

```
Main genome scaffold total: 14932
Main genome contig total:    48818
Main genome scaffold sequence total: 469.0 MB
Main genome contig sequence total:   398.5 MB (-> 15.0% gap)
Main genome scaffold N/L50: 336/323.7 KB
Main genome contig N/L50:    7908/13.8 KB
Number of scaffolds > 50 KB: 1562
% main genome in scaffolds > 50 KB: 83.2%
```
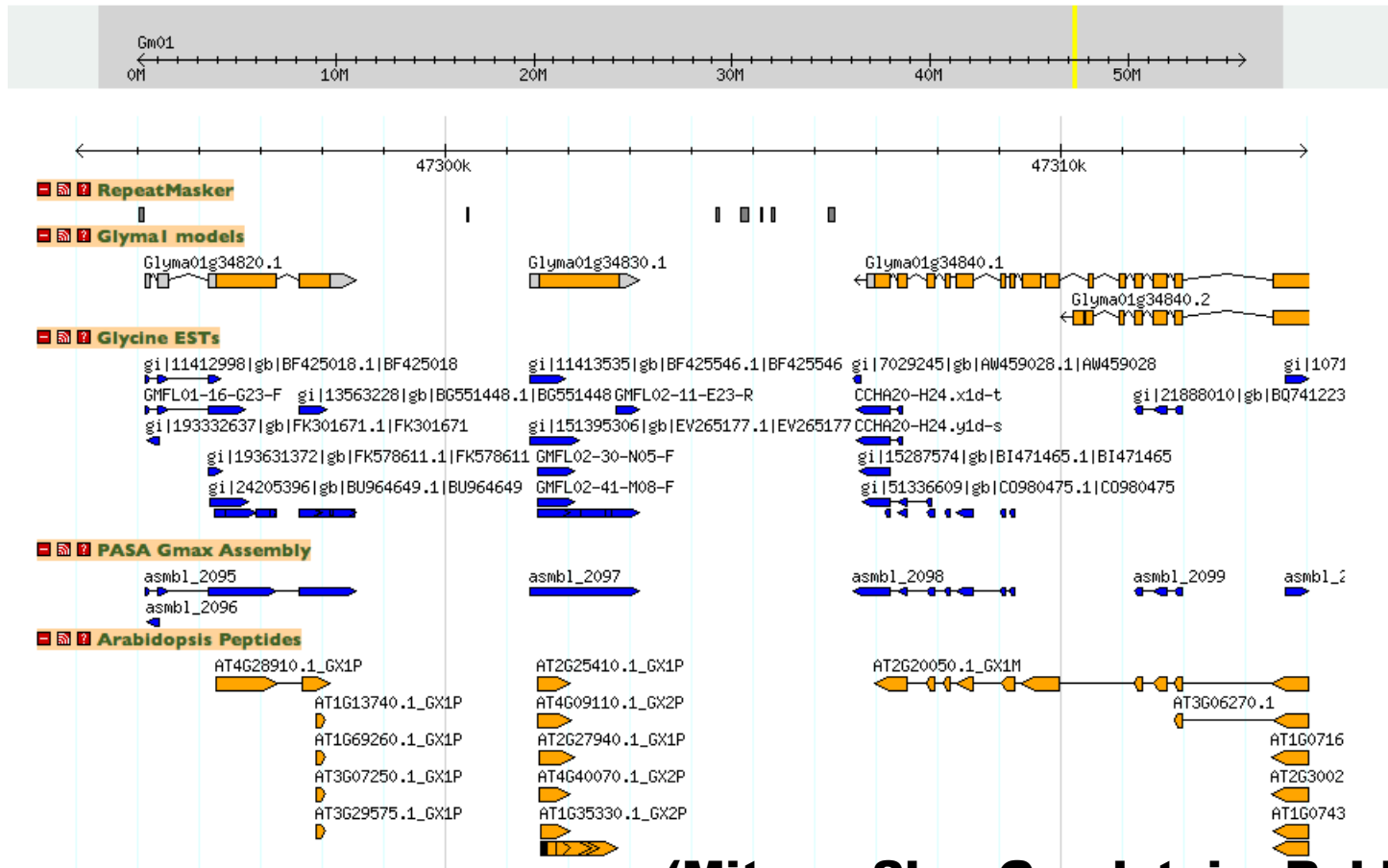
| Minimum Scaffold Length | Number of Scaffolds | Number of Contigs | Total Scaffold Length | Total Contig Length | Scaffold Contig Coverage |
|---|---|---|---|---|---|
| All | 14,932 | 48,818 | 469,034,573 | 398,491,214 | 84.96% |
| 1 kb | 14,932 | 48,818 | 469,034,573 | 398,491,214 | 84.96% |
| 2.5 kb | 10,284 | 44,170 | 458,708,951 | 388,165,740 | 84.62% |
| 5 kb | 5,008 | 38,690 | 441,366,896 | 371,170,537 | 84.10% |
| 10 kb | 3,285 | 35,737 | 430,037,215 | 361,639,508 | 84.09% |
| 25 kb | 2,189 | 32,594 | 412,689,007 | 347,700,588 | 84.25% |
| 50 kb | 1,562 | 29,642 | 390,023,800 | 330,043,578 | 84.62% |
| 100 kb | 951 | 25,234 | 346,755,218 | 296,695,458 | 85.56% |
| 250 kb | 450 | 18,194 | 267,049,160 | 232,736,465 | 87.15% |
| 500 kb | 185 | 11,205 | 175,356,909 | 155,374,306 | 88.60% |
| 1 mb | 56 | 5,205 | 86,489,208 | 77,731,920 | 89.87% |
| 2.5 mb | 5 | 779 | 14,735,915 | 13,561,325 | 92.03% |
| 5 mb | 0 | 0 | 0 | 0 | 0.00% |

# New assembly pipeline

- Remove duplicate 454 pairs and unpaired reads
- Correct 454 bp and indel errors using Illumina reads
- Assemble using HA Arachne version including BES and FES and all of the 454 corrected sequence data
- Make scaffold breaks on misjoins by comparing to genetic map
- Order and orientate scaffolds based on genetic map and build chromosome scale assembly
- This same pipeline is being used to produce multiple plant genomes for the JGI Plant Genome Program

# Annotation

- PV will be annotated using the same Phytozome pipeline as soybean, our goal is to have a comparable gene annotation between the two genomes



**(Mitros, Shu, Goodstein, Rokhsar)**

# Timeline goals

- April 2011: Data collection complete
- July 2011: Chromosome scale assembly
- September 2011: Annotation complete and V1 PV publicly available