

Applied bioinformatics in genomics

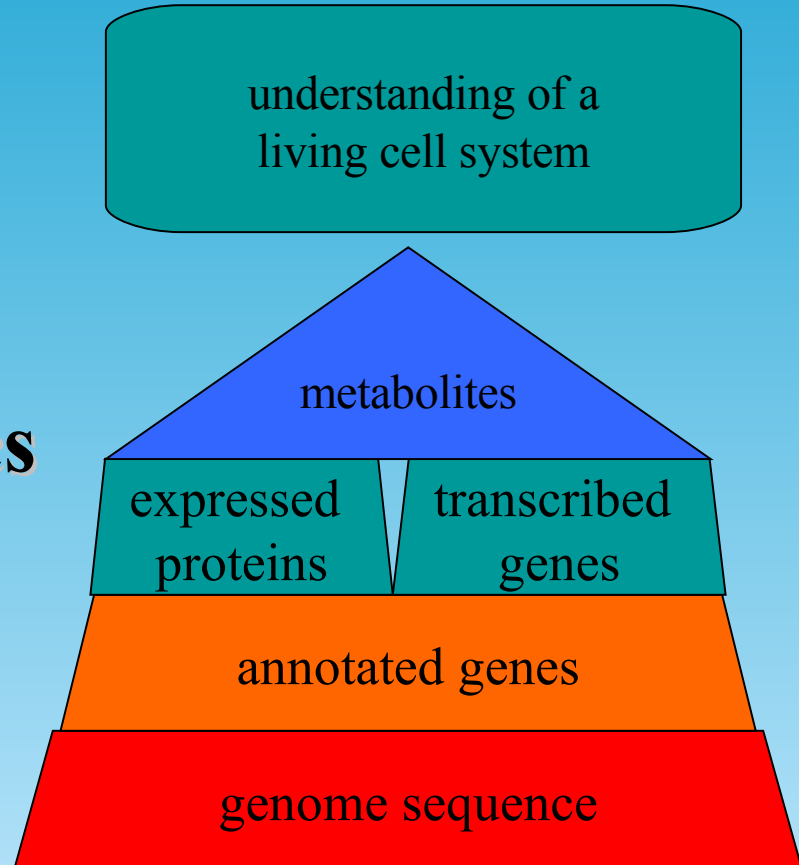
Productive bioinformatics in a
genome sequencing center

Heiko Liesegang

Warschau 2005

The -omics pyramid:

1. **Genome sequencing**
2. **Genome annotation**
3. **Transcriptomics**
4. **Proteomics**
5. **Metabolomics**

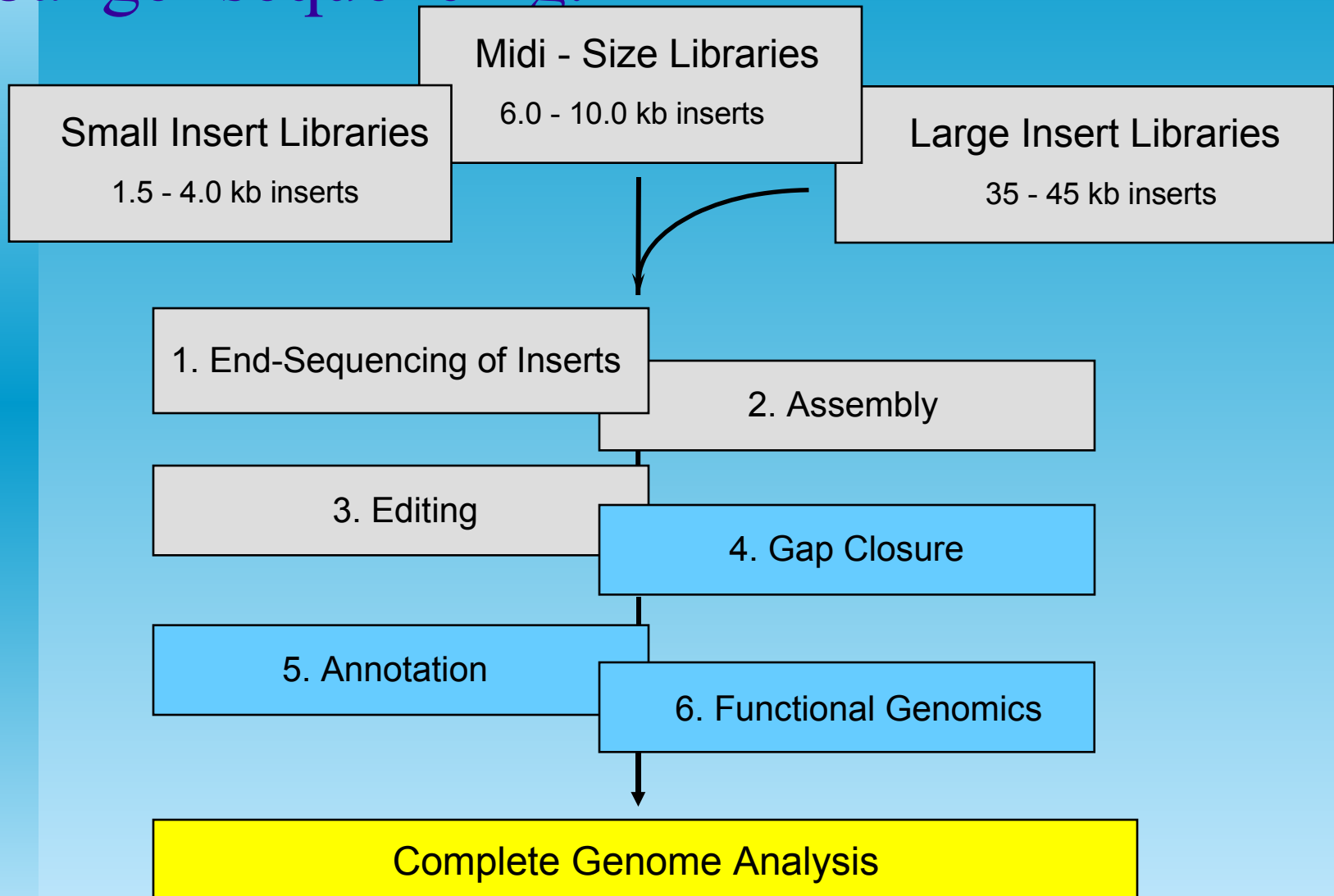


Genome sequencing

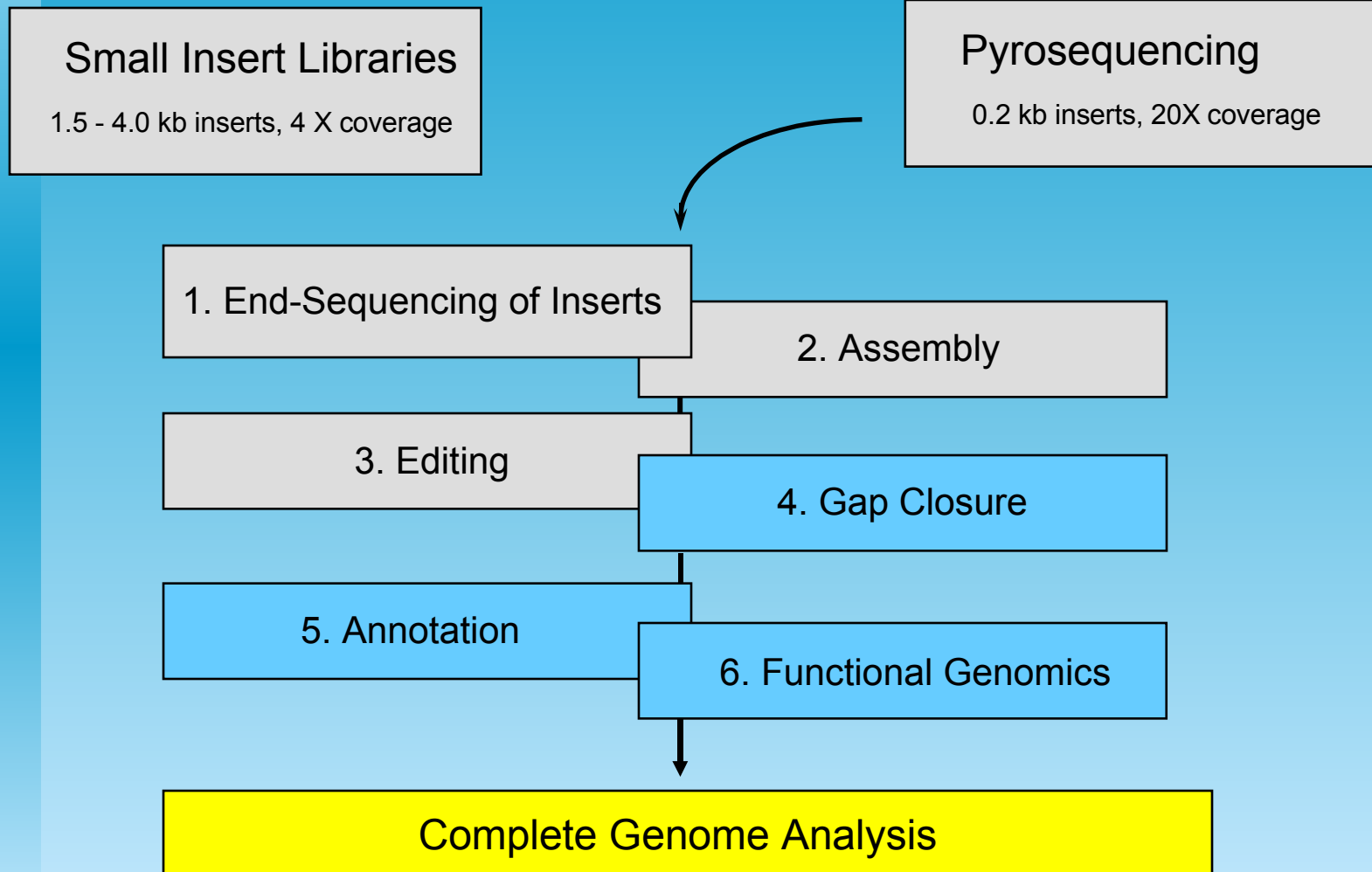
Sequencing Capacities at the G2L

Unit(s)	sequences / week	average read lengths	
1x 454Flex (pyrosequencing)	0,5 x 200.000	~200 bp	60 Mbp (low quality, short reads)
2x ABI 3730XL (96 cap.)	64 x 96	1000 bp	6,14 Mbp (high quality long reads)
	≈ week	≈ Σ 6 Mbp – 66 Mbp	

Genome-Sequencing-Strategy based on Sanger sequencing:

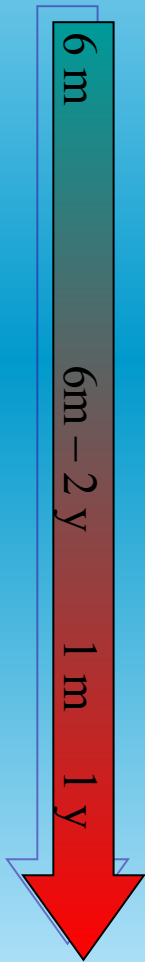


Genome-Sequencing based on mixed Pyro- and Sanger sequencing



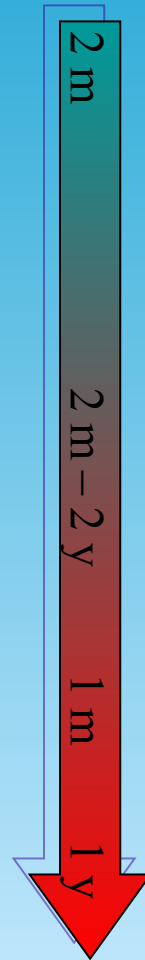
Principle steps in a genome project

Sanger



- **Raw sequencing** – producing random distributed sequences of 8 to 10 fold coverage
- **Processing** – checking the quality of the raw data
- **Assembly** – assemble overlapping reads to reconstruct the replicons of the genome
- **GAP-Closure** – using PCR techniques based on genomic information to generate missing sequences
- **Editing** – manual correction of the computer based assembly
- **Gene prediction** – ORF finding and functional annotation of the DNA-sequence.
- **Annotation** – assign functions to genes to produce biological data

Sanger + 454



- **Raw sequencing A** – producing random distributed sequences of 4 to 5 fold coverage Sanger sequence plus high coverage 454 data ~20 fold coverage
- **Processing** – checking the quality of the sanger raw data
- **Assembly I** – co assemble overlapping reads of sanger sequence with pre-assembled 454 data
- **Assembly II** – reconstruct the replicons of the genome
- **GAP-Closure** – using PCR techniques based on genomic information to generate missing sequences
- **Editing** – manual correction of the computer based assembly
- **Gene prediction** – ORF finding and functional annotation of the DNA-sequence.
- **Annotation** – assign functions to genes to produce biological data

Bioinformatics: Hardware aspects



Requirements of genome assembly:

8 MBp genome assembly with appr. 80,000 reads sanger sequence takes 12 GByte RAM and appr. 2 h calculation time on 1GHz RISC chip

Requirement of genome annotation:

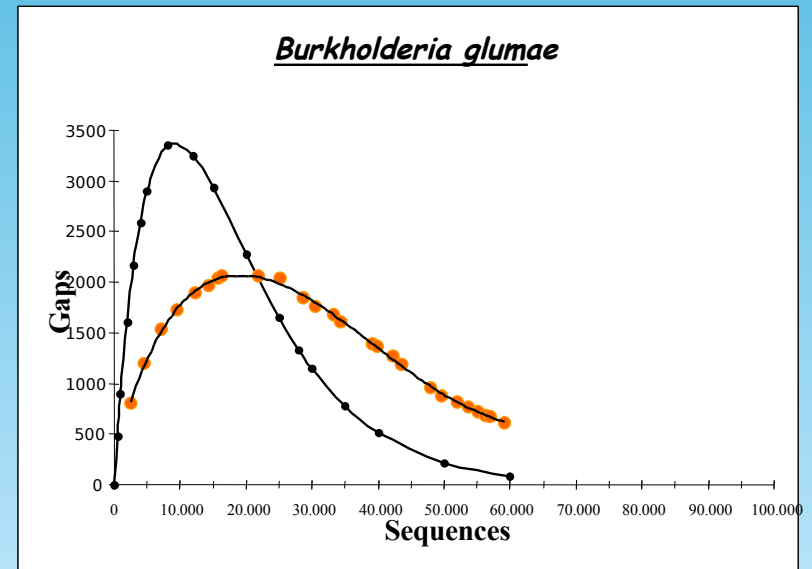
automated similarity based annotation of 1 Mbp genome sequence (~1,000 genes) needs 10 h calculation times on 16 parallel working processors.

Manual annotation:

One human experts annotates 50 genes/d on a web based annotation server

Rawsequencing – the first dataprocessing

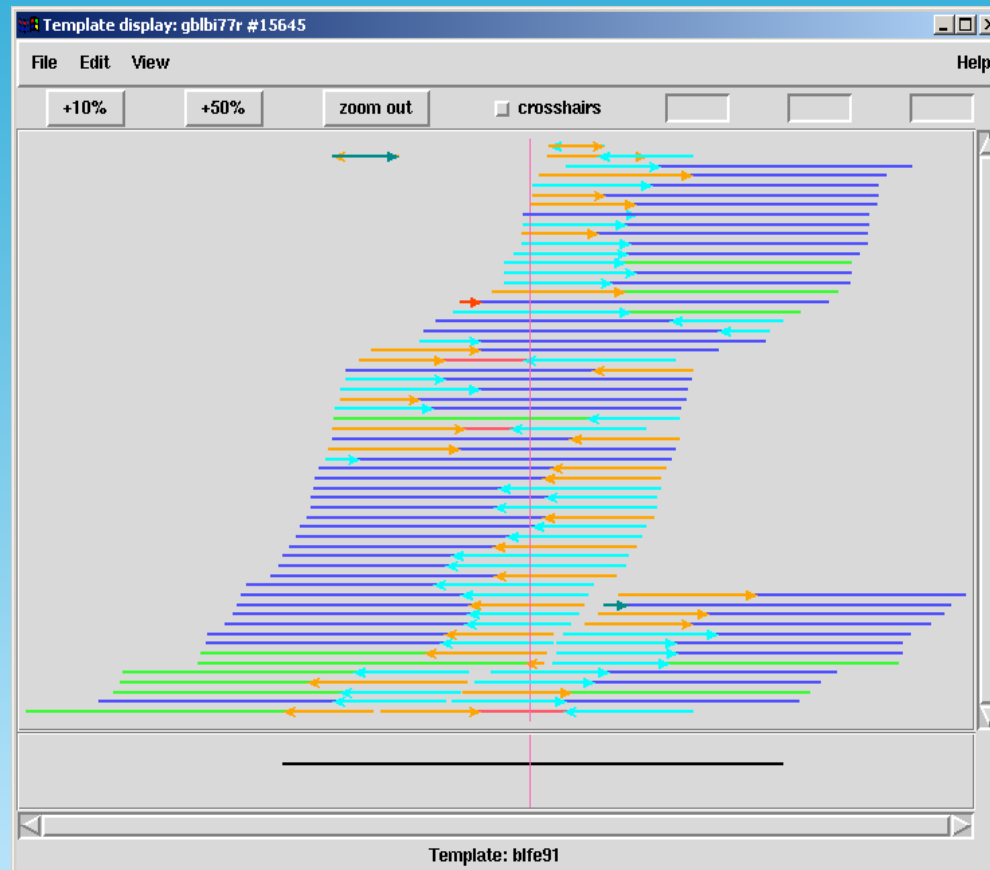
- Check for read quality
- Check gene libraries
- Check genome coverage
- Lander-Waterman plot



Pregap – Screening for contaminating sequences

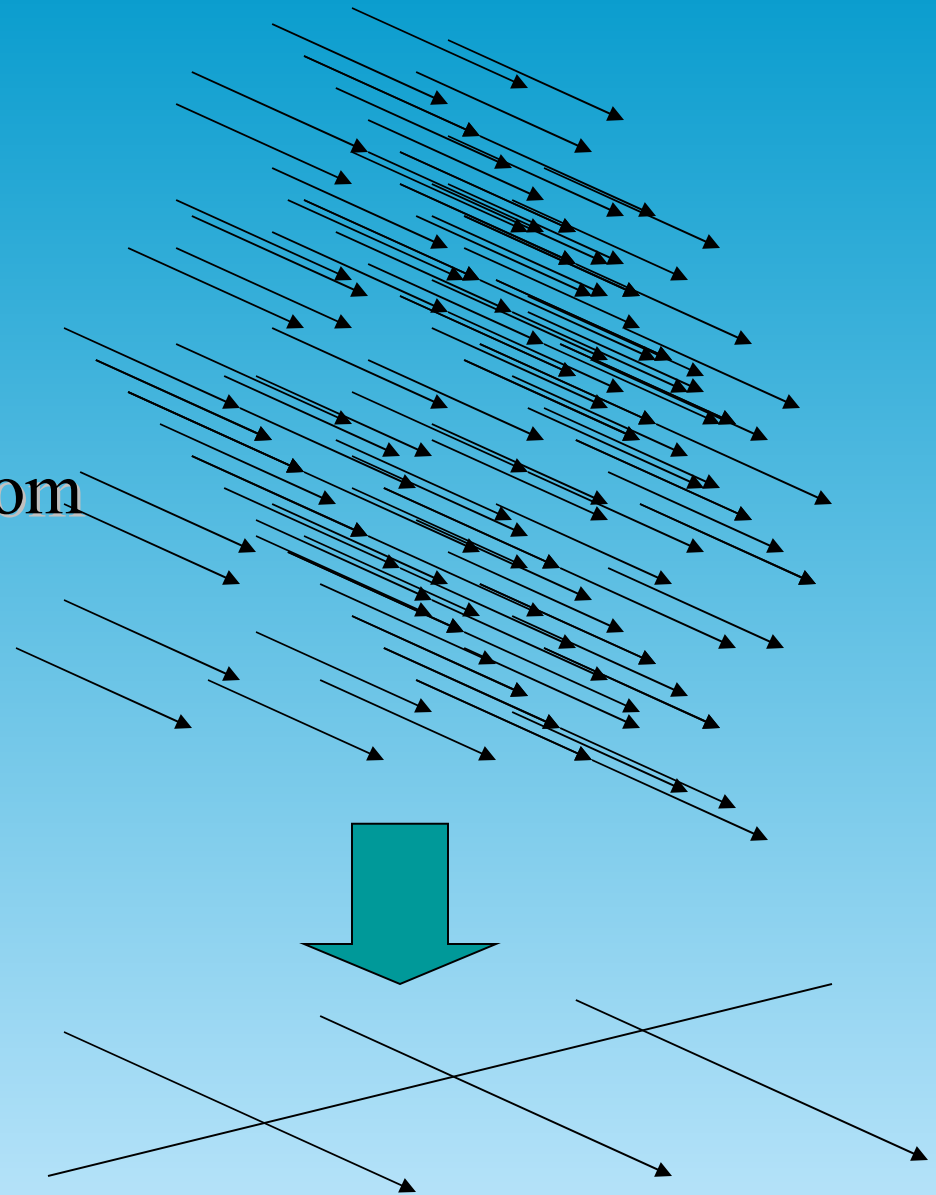
- **Qclip** – defines the minimum quality for segments of sequences used in assembly.
- **Svec** - localises the insertion site at the left end of the sequence.
- **Screenveq** – looks for contaminating vector sequences.

Screenvec failure: sequencing the vector due to gene library contamination



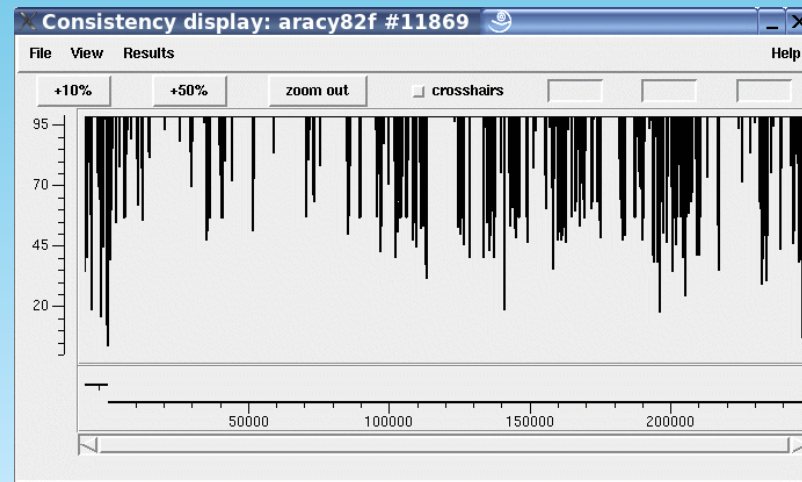
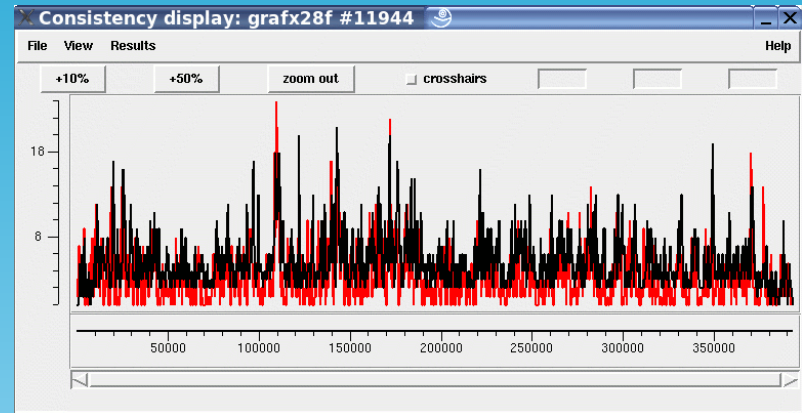
Assembly

- Gcphrap
- Producing contigs from overlapping reads
- Typically 10.000 reads / 1Mbp

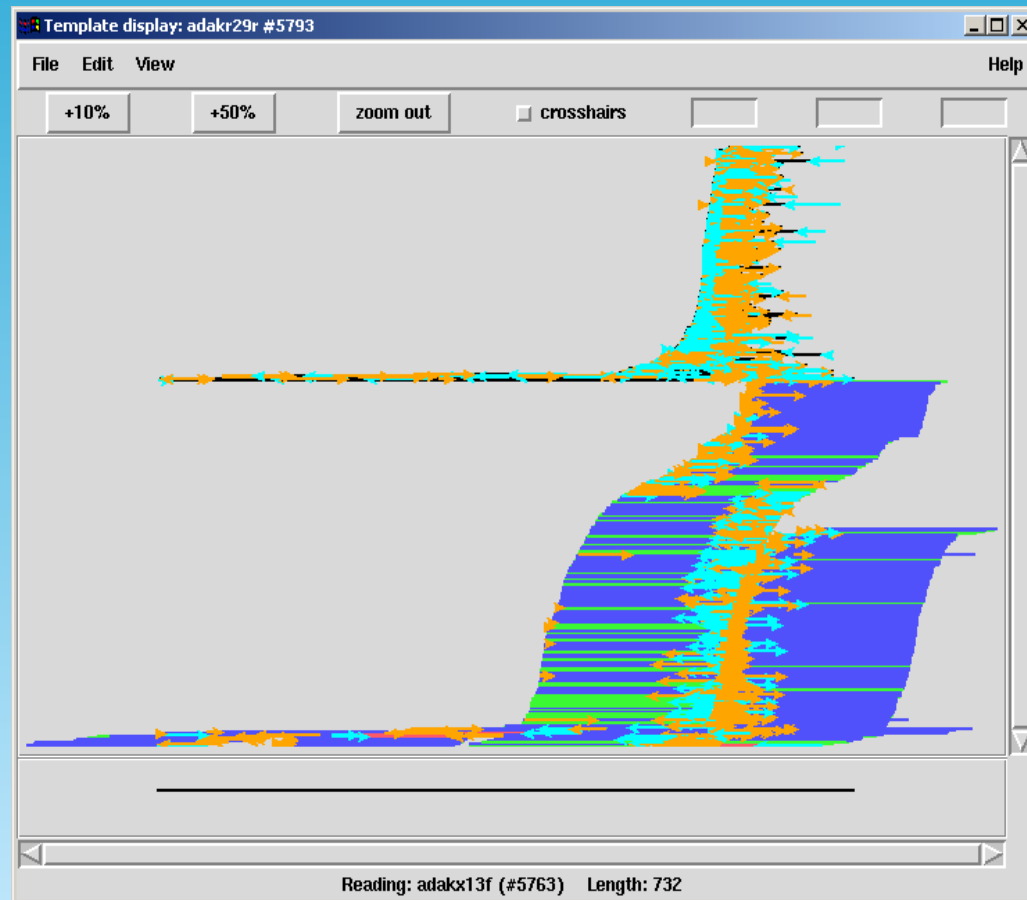


Checking the initial assembly

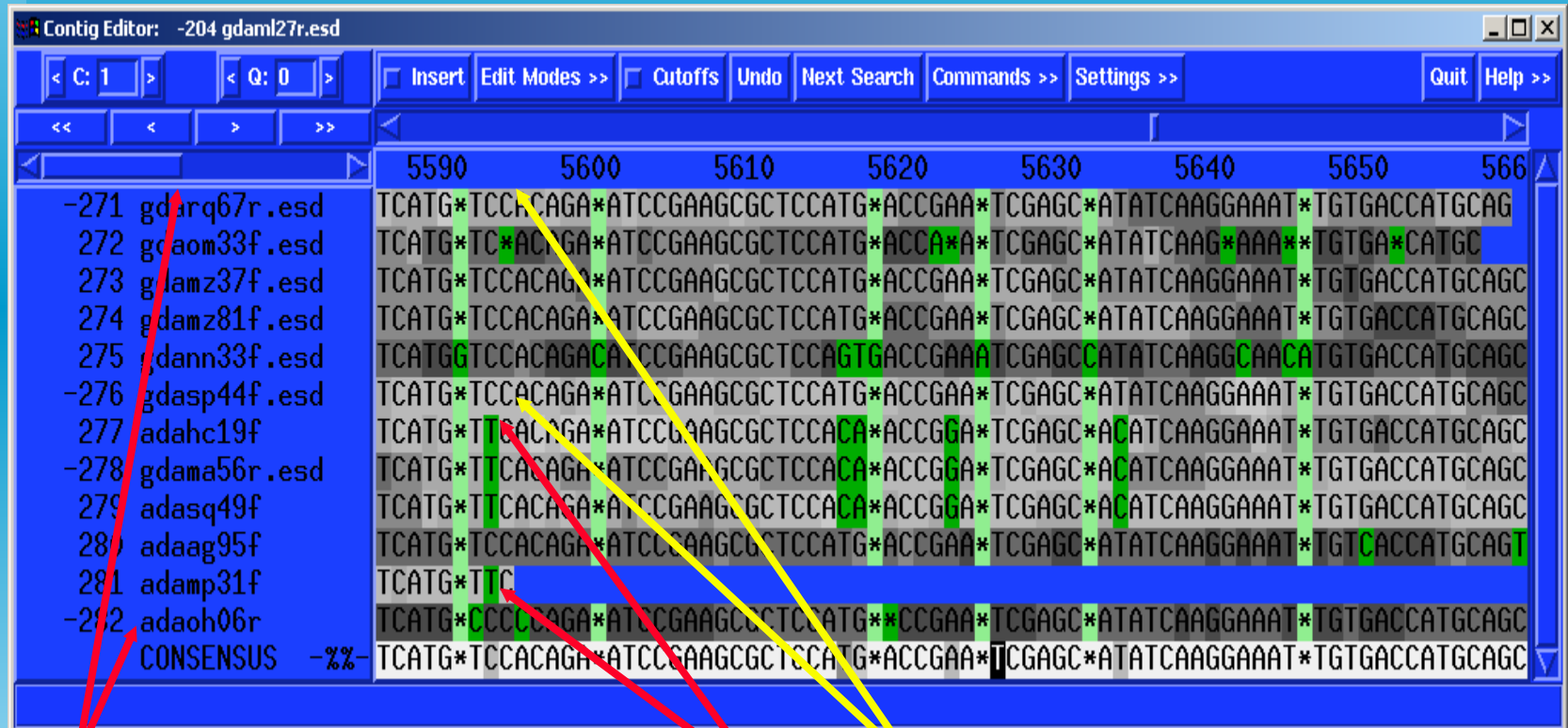
- Distribution of templates
- Size of templates
- Confidence of the contigs



Screenvec failure: Short vector contaminations of $\sim 0.1 - 0.9$ kbp



Missassembly

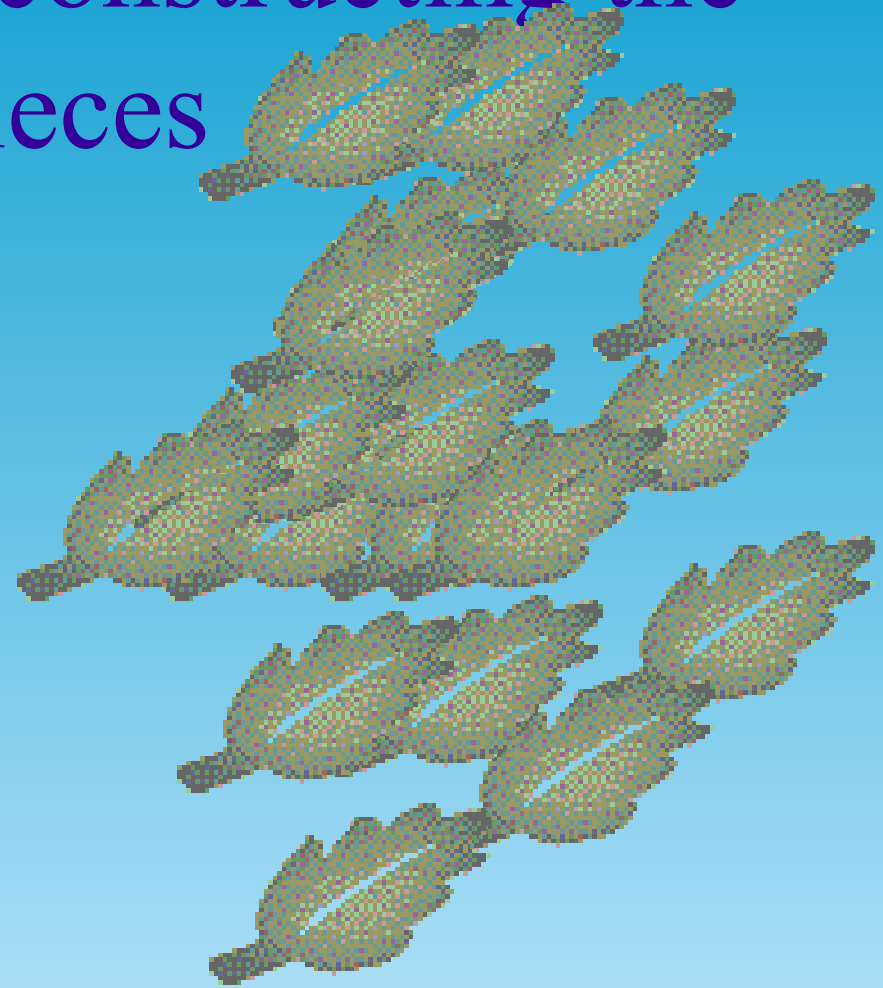


doubled average
base coverage

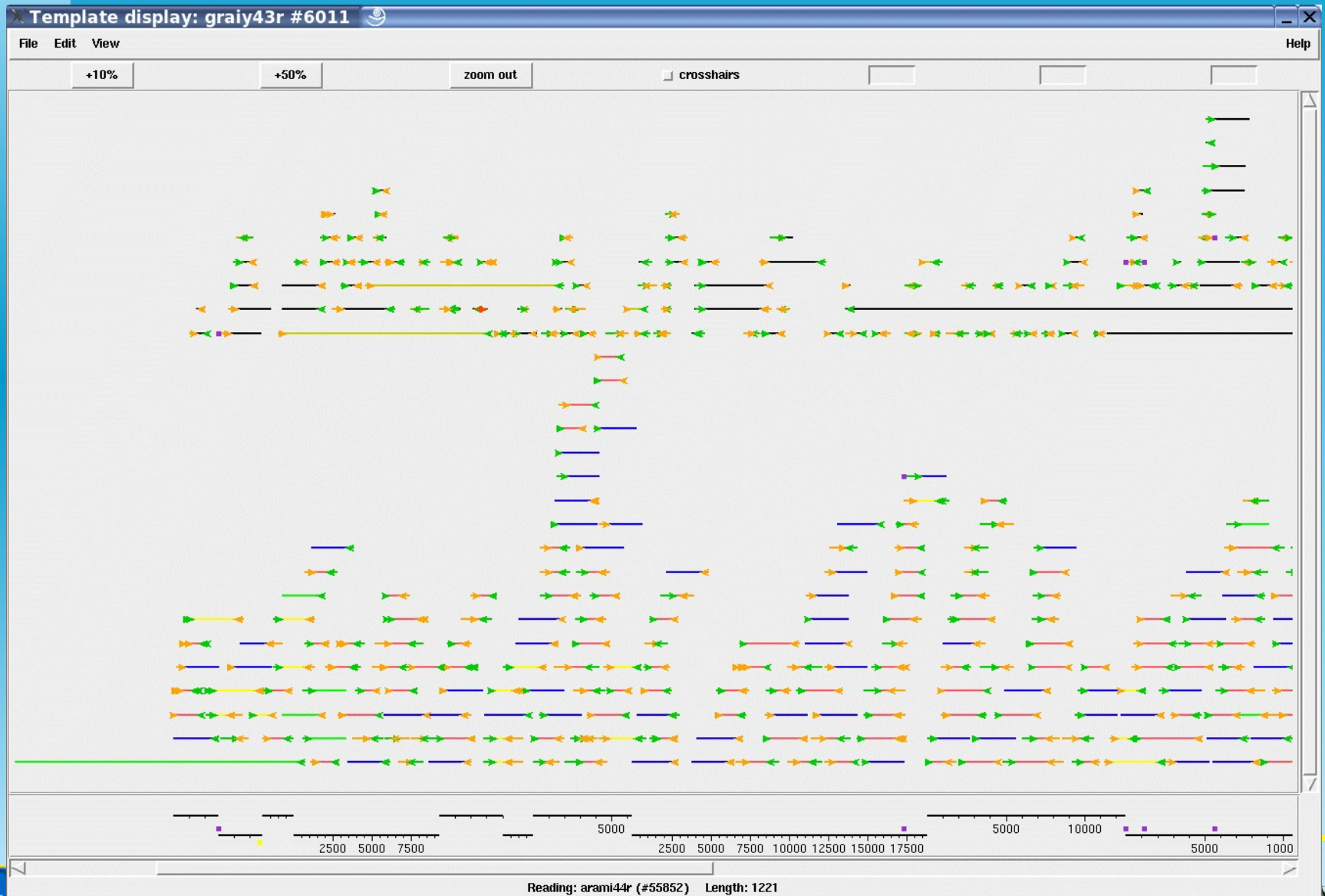
Apparently two groups of
sequences in the multi alignment

GAP-closure reconstructing the genome from pieces

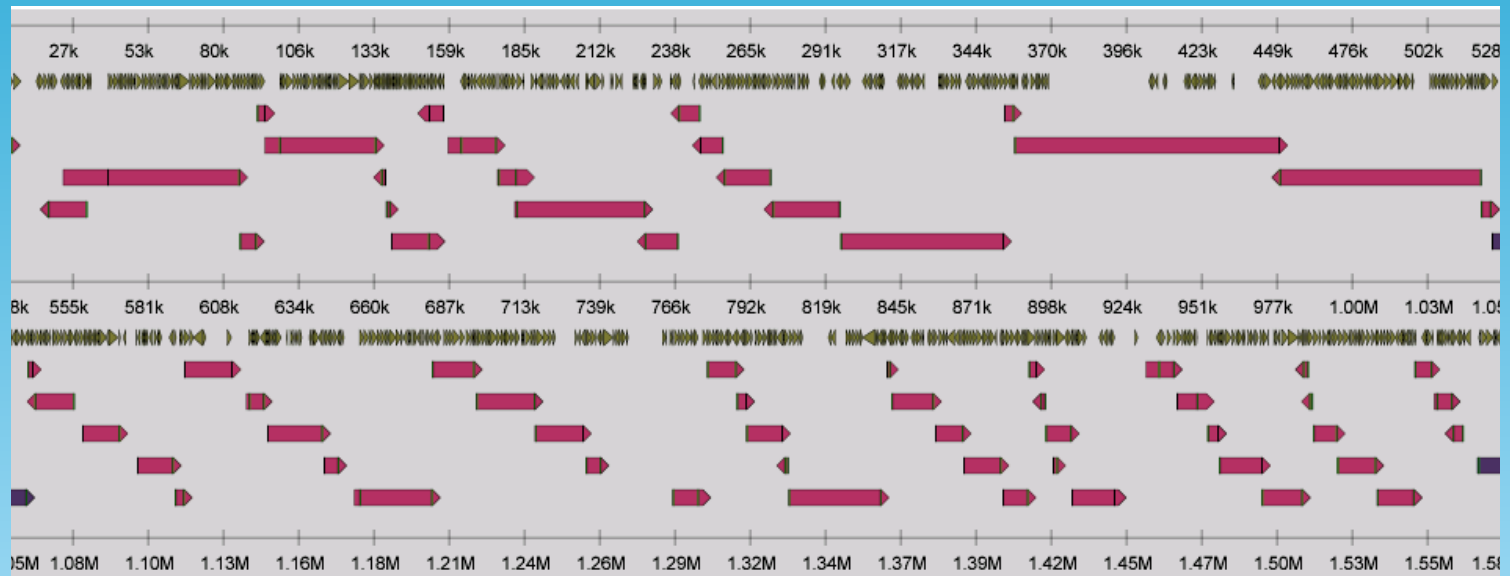
- Sorting contigs
 - ◆ Gap spanning templates
 - ◆ Projector
 - ◆ ERGO gene cluster
- Closing sorted contig gaps by PCR



Sorting contigs by templates



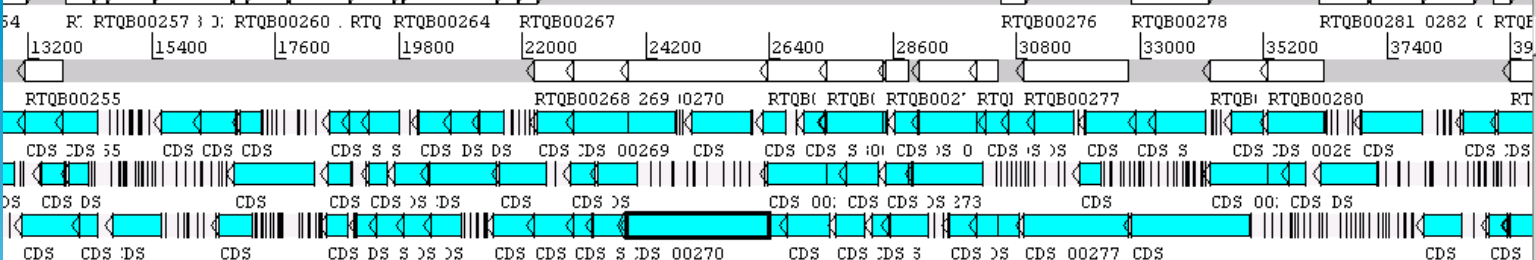
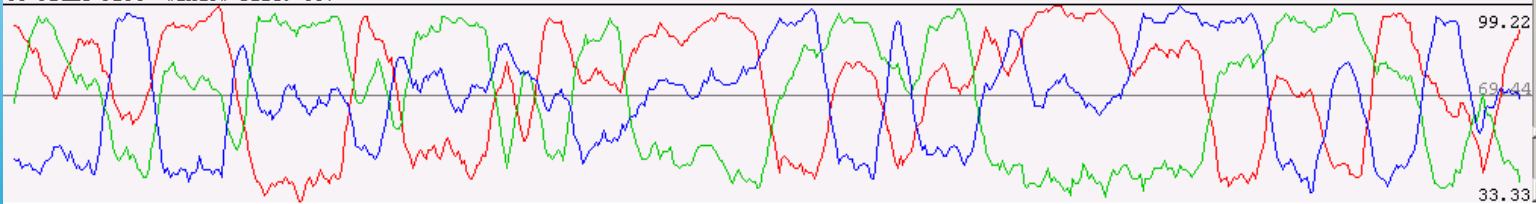
Projector contig projection on related genomes



Selected feature: bases 2472 amino acids 824 RTQB00270 (/gene="RTQB00270" /product="hypothetical conserved protein")

Entry: tbl.chromosome.annotated.gbk ORFS_100+

GC Frame Plot Window size: 387

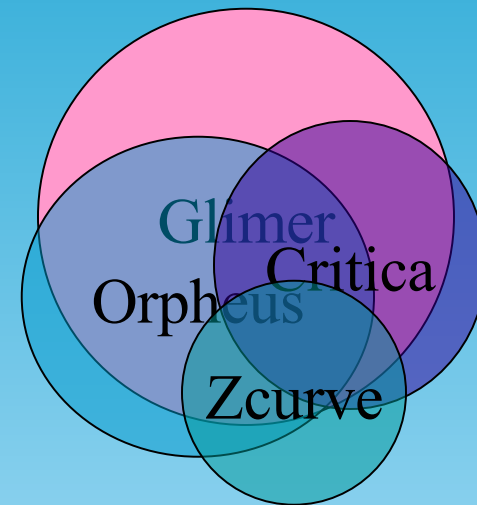
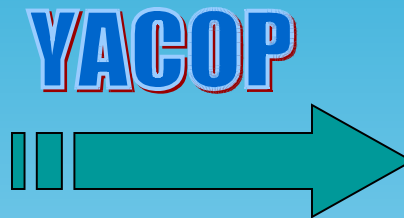


S R T L I T L S S A R A G R R R K R R R G S H L T A A L P G R G K G P R P R R G G W G R R R R G
 P G P * S P S P L R G R E G G G K G E G E A T S P P P S Q E E E K A Q G P V E E V G A E G G E E
 Q D L D H P L L C E G G K E E E K E K G K P P H R R P P R K R K R P K A P + R R L G P K A A R K
 CCAGGACCTTGATCACCTCTCTCTCGGAGGGCGGGAAGGAGGAGGAAAAGGAGAAGGGGAAGCCACCTCACCGCCGCCCTCCAGGAAGAGGAAAAAGGCCAAGGCCCGTAGAGGAGGTTGGGGCCGAAGGGCGGCGAGGAA
 26300 26320 26340 26360 26380 26400 26420
 GGTCCTGGAACTAGTGGAGAGGAGACGCTCCCGCCCTCTCTCTCTTTCTCTCTTCCGTTGGAGTGGCGGGGGAGGGTCTCTCTCTTTCCGGGTTCCGGGGCATCTCTCCAACCCGGCTTCCGCGCTCTCTT
 G P G Q D G E G R R P R S P P P P P S P S A V E G G G E W S S S S F A W P G T S S T P A S S P P S S
 W S R S * G R R Q S P P F S S S F S F P F G G (* R R G G L F L F L G L A G Y L L N P G F A A L F
 L V K I V R E E A L A P L L L L F L L L L P L W R V A A R G P L P F P G L G R L P P Q P R L R R P L

CDS 23895 26366 c
 23895 26366 c

YACOP - Combined ORF finding

- Glimmer
- Orpheus
- Critica
- Zcurve
- ORF finding is based upon statistical properties of coding regions
- Experts have to curate ORFs according to frameshifts, startcodons, artifact ORFfinding



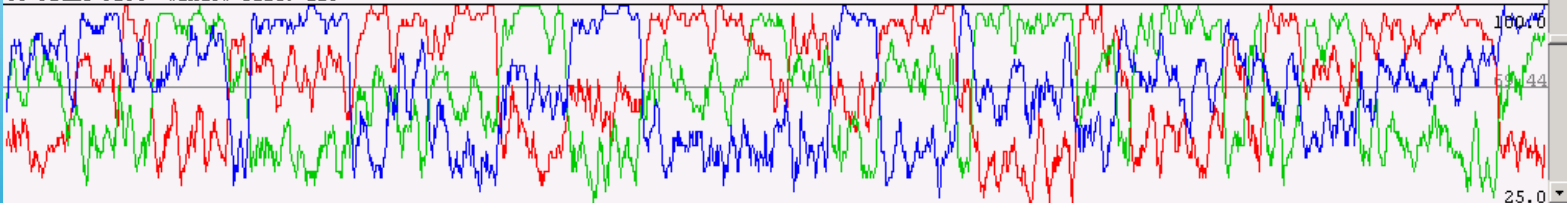
ORF to gene

- *Thermus thermophilus* sequences bare 2 – 3 fold coverage of ORFs per sequence
- Automated ORF finding uses codon biases und HMM to identify potential “coding” ORFs
- Gene start prediction is curated by RBS-finder
- Accuracy: ca. 97% relevant ORFs can be identified
ca. 50% to 65% of gene starts can be correctly predicted
- All ORFs have to be curated by and human experts.

Nothing selected

Entry: tbl.chromosome.annotated.gbk

GC Frame Plot Window size: 120



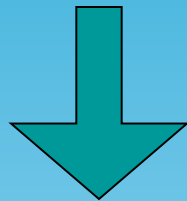
```

L R A A Q K G R Y Q G G M A F V K G P P L * A L I T F L P P V P E A T * S A P P P P L P S K A G
S G Q R K R E D T K E A W P L S R A P L S E L S S L S S P Q F L K Q P D Q L L H H R F P P R R G
. Q G S A K G K I P R R H G L C Q G P P S L S S H H F P P P S S * S N L I S S S T T A S L Q G G
CTCAGGGCAGCGCAAAAAGGGAAGATACCAAGGAGGCATGGCCTTTGTCAAGGGCCCCCTCTCTGAGCTCTCATCACTTTCTCCCCCAGTCTCTGAAGCAACCTGATCAGCTCTCCACCACCGCTTCCCTCCAAGGGGGG
      20          40          60          80          100          120          140
GAGTCCCGTCGCGTTTTCCCTTCTATGTTCTCCGTACCGGAACAGTTCCTCCGGGGGAGAGACTCCGAGTACTCAAGGAGGGGGTCAAGGACTTCGTTGGACTAGTCGAGGAGGTGGTGGCGAAGGGAGGTTCCGCCCC
E P C R L L S S V L S A H G K D L A G R E S S E D S E E G W N R F C G S * S R W W R K G G L R P
* P L A F P F I G L L C P R Q * P G G E R L E * * (K G G G L E Q L L R I L E E V V A E R W P P P
. L A A C F P L Y W P P M A K T L P G G R Q A R M V K R G G T G S A V Q D A G G G G S G E L A P
    
```

source 1 1894877
cns 22 1121

T. thermophilus

- *T. thermophilus* contains two replicons :
 - ◆ Chromosome (1,894,877 bp, 69.2% GC content)
 - ◆ pTT27(232,605 bp, 69.4% GC content)

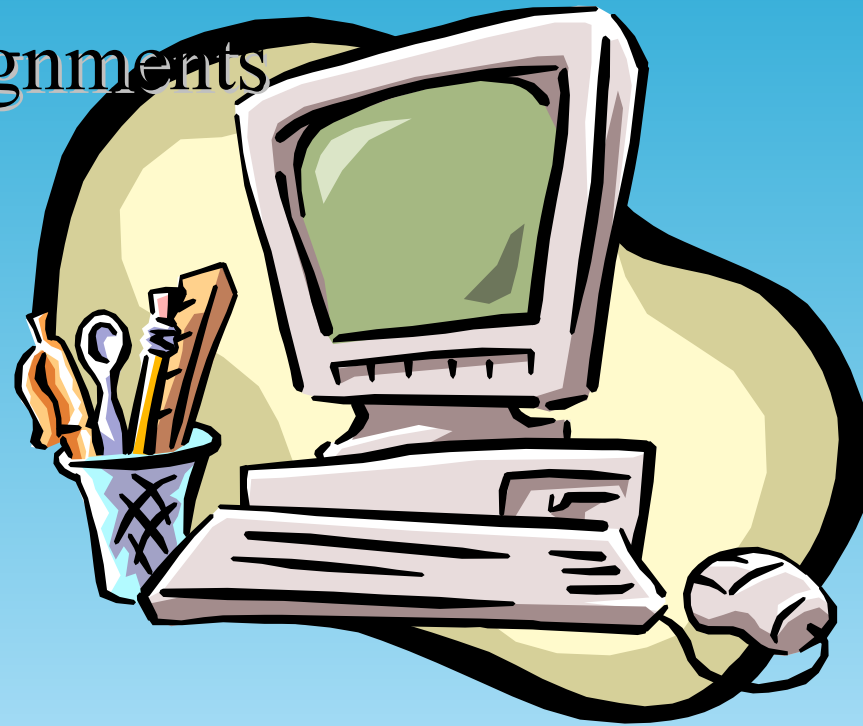


Orffinding,
manual curation

ORFs	Chromosome	pTT27
with assigned function	1397 (70%)	85 (37%)
conserved hypothetical	192 (10%)	56 (24%)
no database match	399 (20%)	89 (39%)
Total	1988	230
rRNA – Cluster	2	-
tRNA	47	-
IS-Elements	30	23

The first Annotation steps- similarity based automated assignments

- FASTA - automated assignments
 - ◆ for all proteins against a high quality dataset SwissProt, PIR
- BLAST, PSI-Blast
 - ◆ Providing data for manual control



Transitive error propagation

- In automated annotation system as
- Genequiz
- PEDANT
- ERGO/WIT
- Annotations are assigned if proteins share a cut of similarity
- => original errors thus multiply by each new entry



Manual Curation of automated Annotations

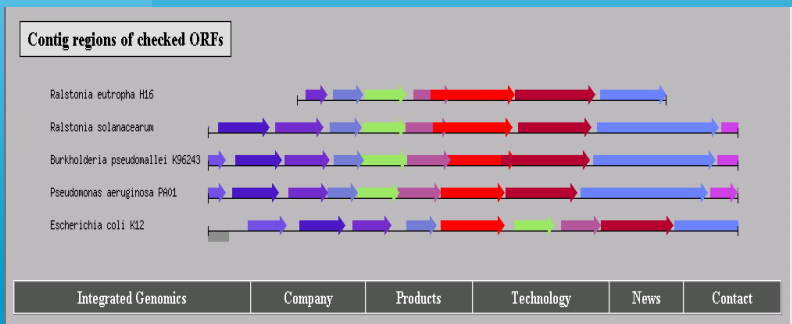
- A team of annotation experts has to check each automated annotation
- ERGO provides additional non similarity based means to evaluate Annotations
- Tmpred, Pfam, COGs, Prosite are used to achieve a consistent functional assignment
- The kind of assignment evidence is a way of classifying annotation reliability

The screenshot displays the ERGO web interface. At the top, the ERGO logo is shown with the text "427 genomes Statistics User hliese Integrated Genomics". Below the logo is a navigation bar with "Data", "Query", "Configure", and "Help" menus. The current page is titled "Protein Page for RBL00098 from Thermus thermophilus". The main content area is divided into two columns. The left column contains a "Data Panel Display" section with a "Select Data Panel" dropdown and an "Examine RBL00098" section with a "View Annotations" link. The right column contains "Primary Information for RBL00098 Thermus thermophilus" and "Contig Region for RBL00098".

Primary Information for RBL00098 Thermus thermophilus	
Aliases	None
Contig Location	chromosome from 1,090,978 to 1,092,039; contig length =
AA Residues, DNA	354 aa, 1062 bp (returns sequence)
Molecular Weight	39,716 (returns all proteins with Molecular Weight \pm 1 per
Iso-electric Point	5.94 (returns all proteins that have pI within \pm 0.10)
Function	Uroporphyrinogen decarboxylase (EC 4.1.1.37)
EC 4.1.1.37 Links	Enzyme Commission, Expaty, KEGG Maps

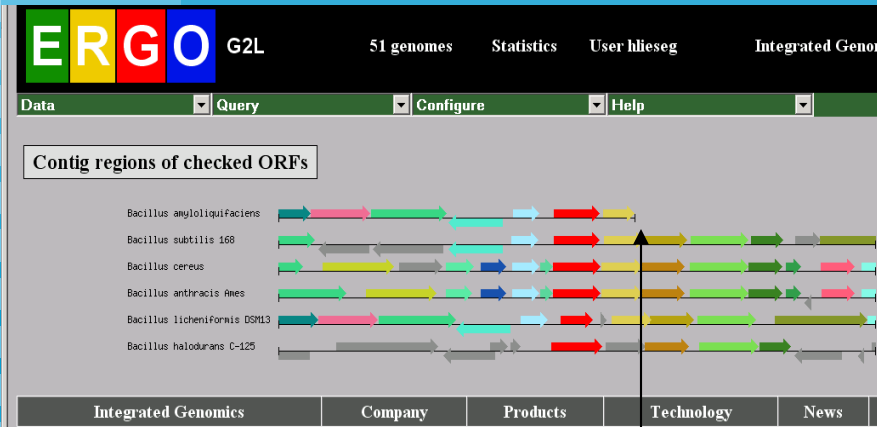
Contig Region for RBL00098	
Neighboring Genes	1,081,508 bp

Functional coupling in genomes

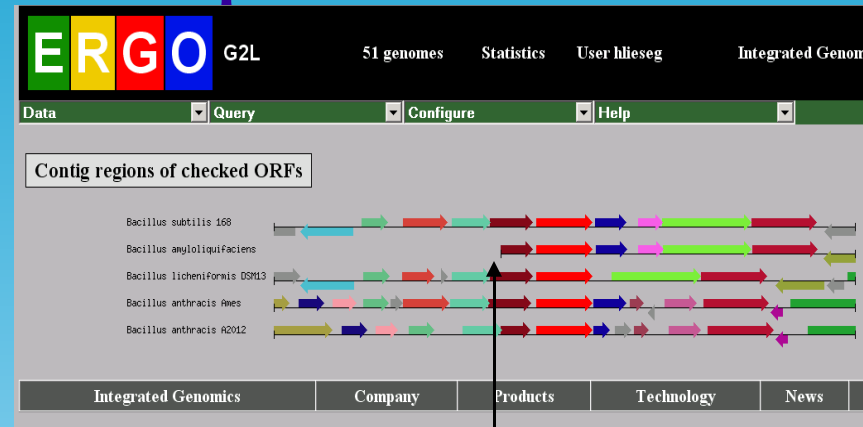


- Most bacterial genes are clustered.
- Related organisms share the organization of gene clusters
- Cluster provide a context of genes for reannotation
- Cluster provide a datasource for gapclosure

Alignment of clustered genes from distant related species



End of one
contig inside a
gene cluster



End of another
contig inside the
same gene cluster

Metabolic reconstruction

- Based on Annotated genes metabolic pathways can be predicted
- Pathway finding is based upon EC categorisation
- Missing EC functions have to be searched in the “hypothetical” genes
- Thus pathway reconstruction leads to a reannotated genome

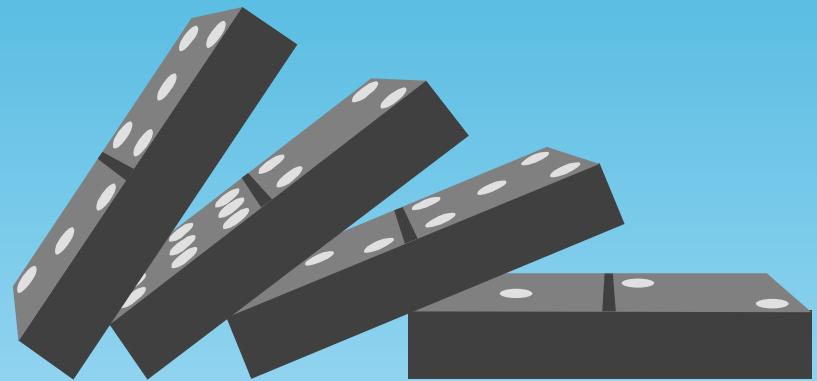
Genomic properties

- Genome sequences provide physiological predictions
- Predicted functions have to be proven experimentally
- High throughput evidence methods from transcriptomics and proteomics give exists/not exists data
- Functional evidence have to be provided by classical wet lab experiments



Source of annotation errors

- ORF-finding errors
- Sequencing errors
- Artifact similarity
- annotation errors in reference database
 - ◆ Misspelling
 - ◆ Underannotation
 - ◆ Overannotation
 - ◆ Missannotation



Error in published genomes: Enzymes an EC classification

- The EC classification for enzymes characterizes Enzyme by a four digit number
- EC 6.1.1.19 (=) Arginyl-tRNA Synthetase
(first characterized as Arginyl-tRNA Transferase, a complete different reaction type)
- **(enzymetype).(cofactors).(reactiontype).(substrate)**
- Errors found:
- First digit class of enzyme (~8 % errors)
- Last digit substrate specificity (\geq 33 % errors)

Scaling of annotation errors: Transitive Annotation-Based Scale

■ Description	Score	description
■ False positive	7	wrong function
■ Over-prediction	6	over detailed annotation
■ Domain error	5	domain error
■ False negative	4	function not found
■ Under prediction	3	under detailed annotation
■ Undefined	2	-
■ Typographical error	1	misspelling

Chlamydia trachomatis in different annotations

Score	Original	Genequiz
7	43	35
6	50	62
5	20	23
4	26	50
3	67	48
2	84	64
1	35	10
0	565	598

Ranking of reliability

1. Biochemical experimental characterization
2. Genetic experimental data from mutants
3. Consistent properties to a highly similar database hit
4. Good database hit
5. Consistent frame according to the used ORFfinder
6. Nonoverlapping ORF of minimum size

YACOP – gene prediction pipeline in G2L

- Gene prediction is based upon different aspects of typical genes.
- YACOP uses an combined ORF-finding of Glimmer, Zcurve and Critica with a startcodon correction by RBS-finder
- All ORFs a curated manually by human experts.
- In genetic islands genes have still to be predicted by human experts, automated ORF-finding fails almost complete.

Gene prediction – looking for biological features

Protein coding ORFs ca. 90 – 99 % (only 65% correct starts)	YACOP, TICO, GeneMark, Glimmer, ...
tRNA ca. 99%	tRNA-Scan
rRNA ca. 99%	BLASTN
ncRNA ?	structure based evaluation of short alignments

Annotation

- Assigning functions to sequences
- **Proteins:**
- GenDB
- ERGO
- **Islands/prophages:**
- SIGI/Colombo
- **Transcriptomics:**
- Which genes are working?