

GTrack 1.0:

**Unified data format providing
customizable representation and
high-speed analysis performance
within Galaxy**

Sveinung Gundersen, PhD student
The Norwegian Radium Hospital, OUH

September 3rd, 1967, Stockholm, Sweden



from <http://retronaut.com>

The day Sweden changed from driving on the left to driving on the right

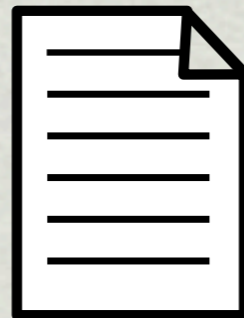
File formats for genomic track data



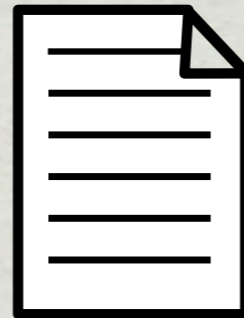
BED



bigBed



SAM



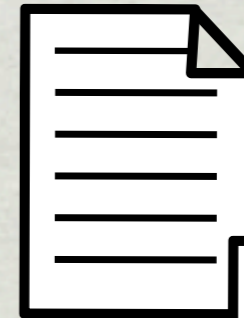
FASTA



GFF



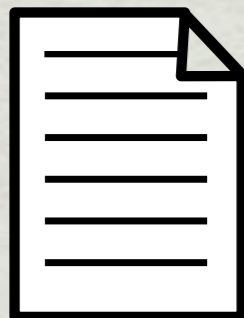
bigWig



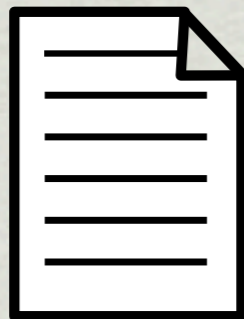
VCF



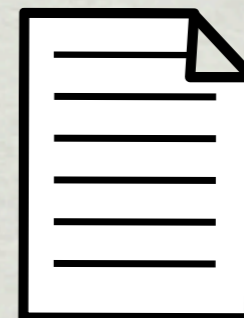
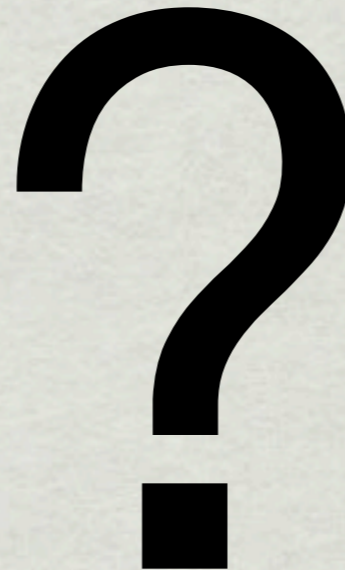
broadPeak



WIG

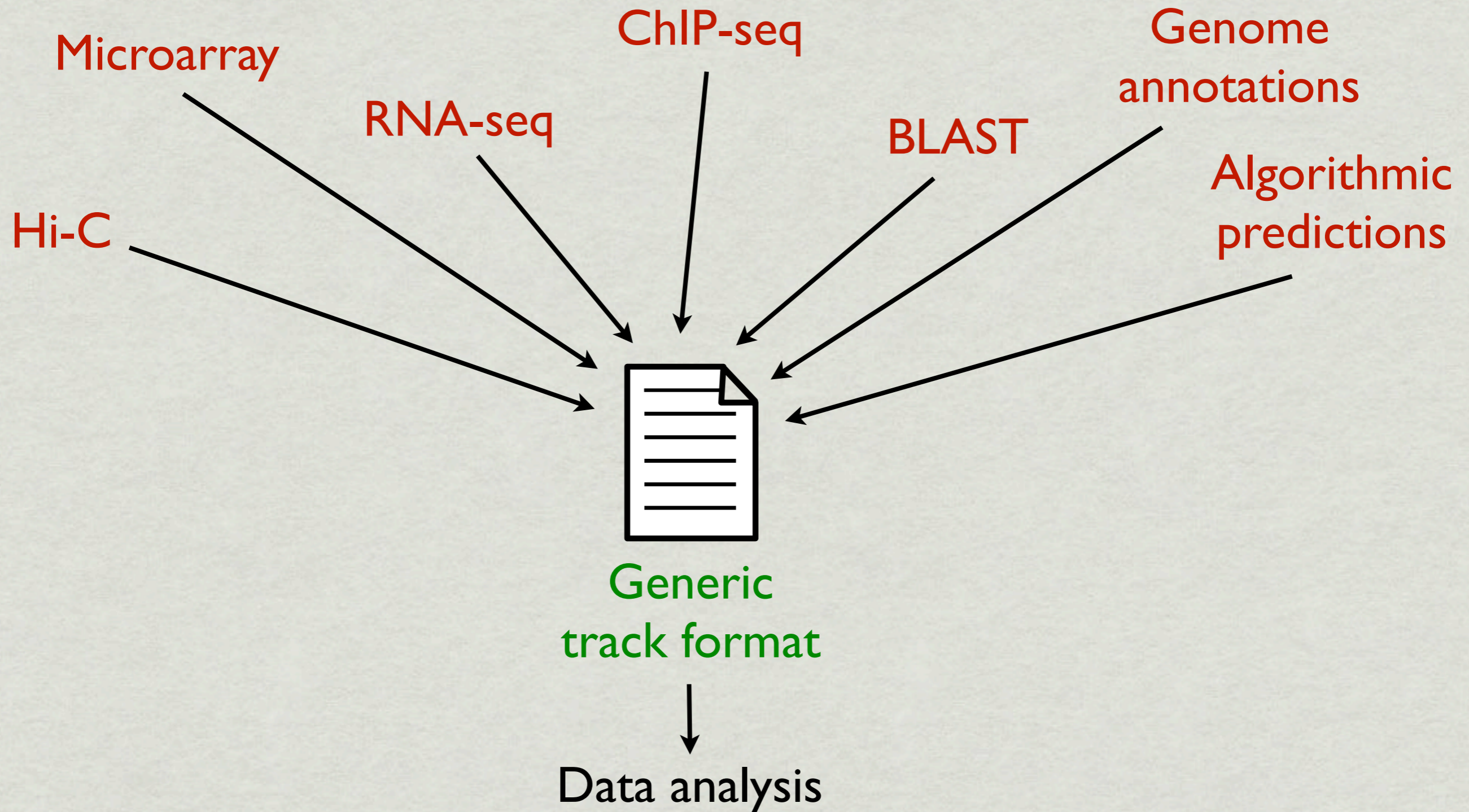


Galaxy
interval



BedGraph

Why couldn't our world just be like this?



Why all the formats?

1. Different columns
2. Underlying structural differences
3. Format-specific characteristics
4. Tool-specific format requirements
5. Performance issues

Why all the formats?

1. Different columns
2. Underlying structural differences
3. Format-specific characteristics
4. Tool-specific format requirements
5. Performance issues

Different columns

Different columns

- Different set of columns in the different formats

Different columns

- Different set of columns in the different formats
- Solution: handle flexible columns in one format

Different columns

- Different set of columns in the different formats
- Solution: handle flexible columns in one format
- Galaxy interval format handles flexible columns:

#CHROM	START	END	STRAND	NAME	COMMENT
chr1	10	100	+	exon	myExon
chrX	1000	10050	-	gene	myGene

Different columns

- Different set of columns in the different formats
- Solution: handle flexible columns in one format
- Galaxy interval format handles flexible columns:

#CHROM	START	END	STRAND	NAME	COMMENT
chr1	10	100	+	exon	myExon
chrX	1000	10050	-	gene	myGene

- However:
 - columns are specified as Galaxy metadata (*i.e.* not readily exportable to other tools)
 - only supports interval data

Different columns

- bigBed also handles flexible columns
- However:
 - columns are specified in server-side AutoSQL files
 - only supports interval data

Why all the formats?

1. Different columns
- 2. Underlying structural differences**
3. Format-specific characteristics
4. Tool-specific format requirements
5. Performance issues

Underlying data type differences

Underlying data type differences

- Example: one value per base pair

Underlying data type differences

- Example: one value per base pair
- BedGraph format:

```
chr1    0    1    1.0  
chr1    1    2    -1.3  
chr1    2    3    2.4  
...
```


Underlying data type differences

- Example: one value per base pair
- BedGraph format:

```
chr1    0    1    1.0  
chr1    1    2    -1.3  
chr1    2    3    2.4  
...
```

- WIG

```
fixedStep chrom=chr1  
1.0  
-1.3  
2.4  
...
```

Underlying data type differences

- Example: one value per base pair

- BedGraph format:

Size (Human genome)

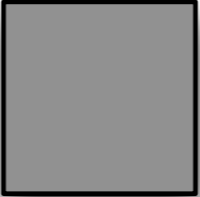
```
chr1    0    1    1.0  
chr1    1    2    -1.3  
chr1    2    3    2.4  
...
```

100 GB

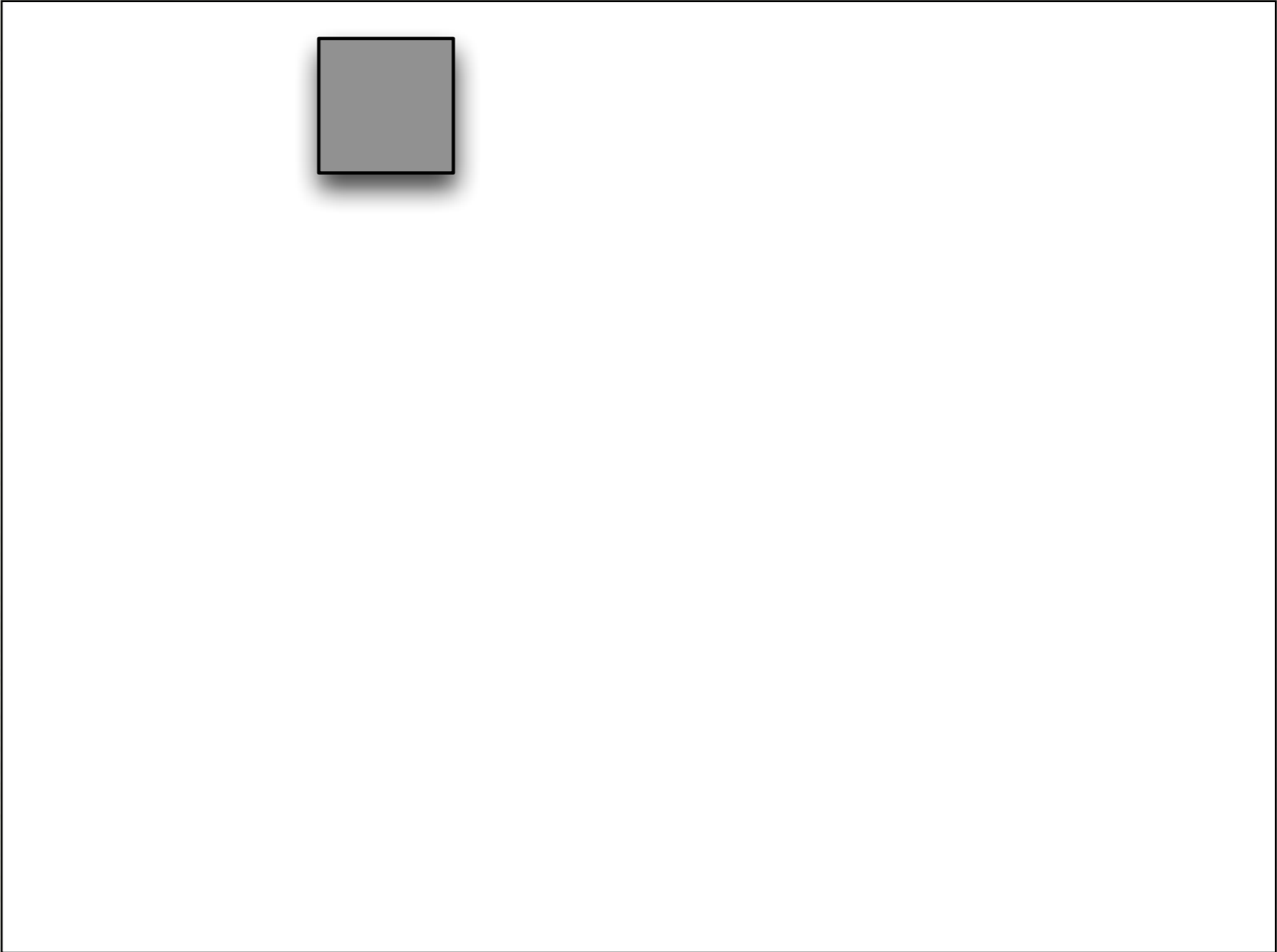
- WIG

```
fixedStep chrom=chr1  
1.0  
-1.3  
2.4  
...
```

20 GB



No
information

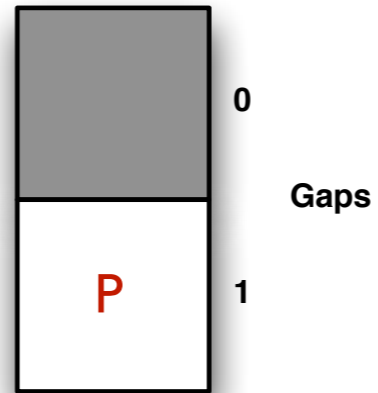


Track
geometry

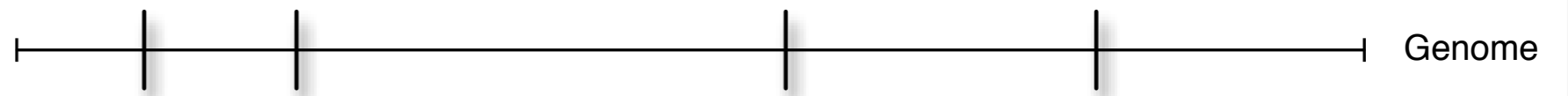


Track
type #1:

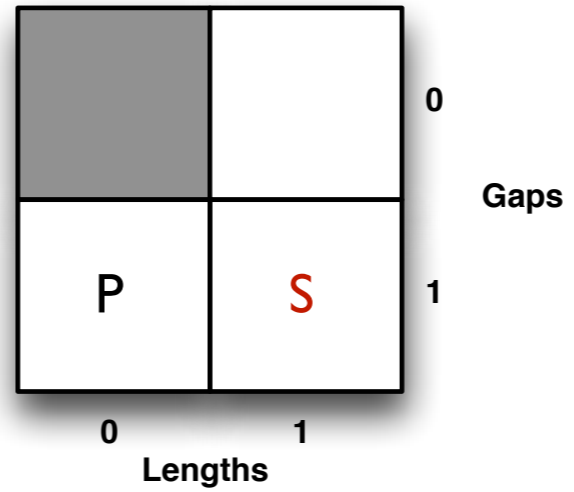
Points



Track
geometry

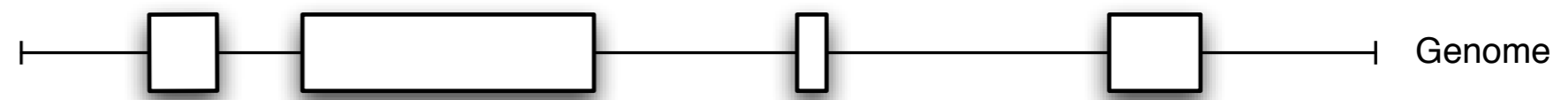


Track
type #2:



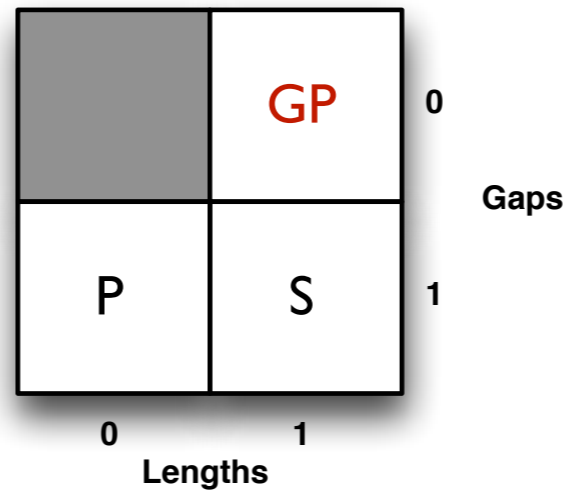
Segments

Track
geometry

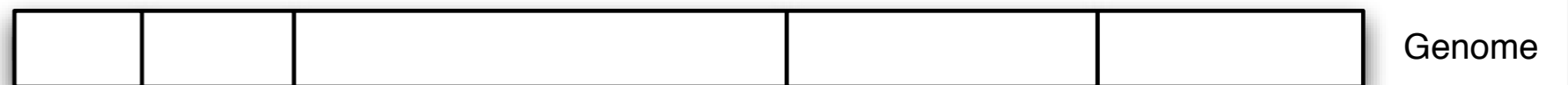


Track
type #3:

Genome
partition

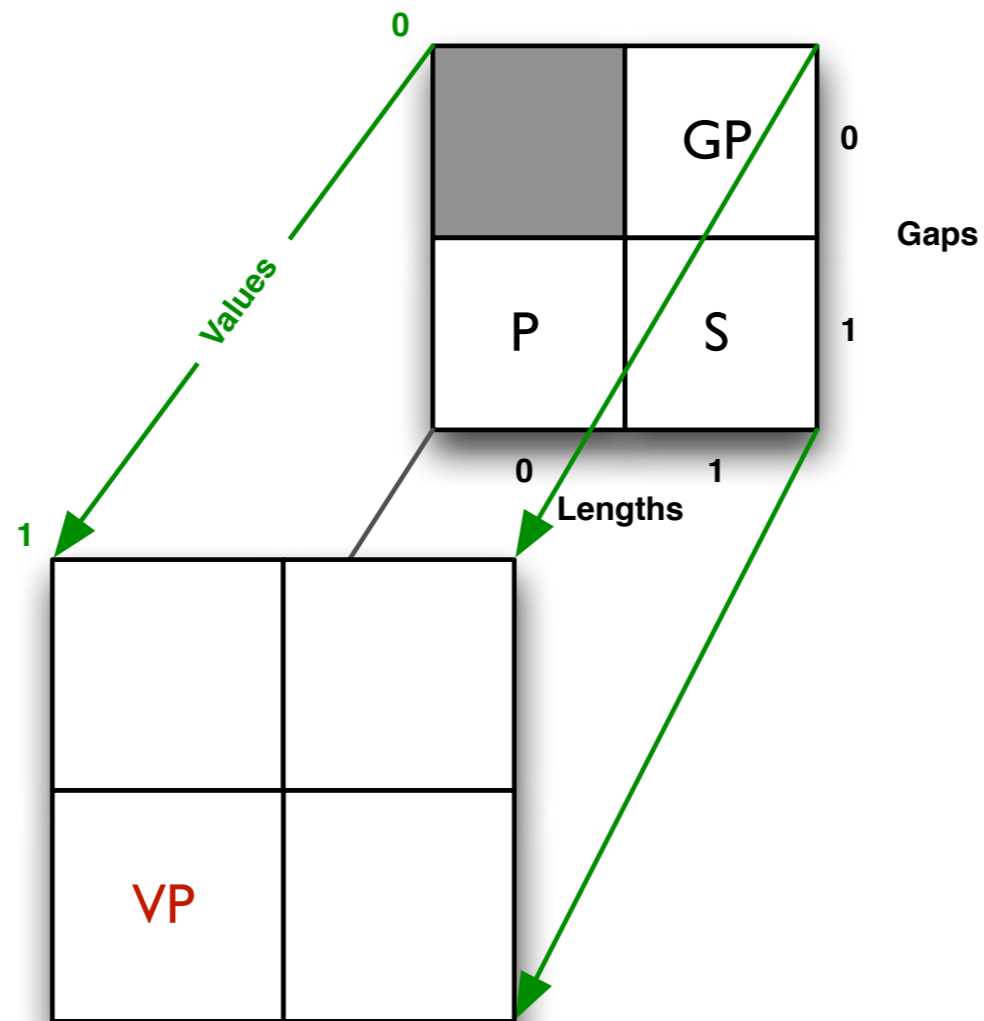


Track
geometry

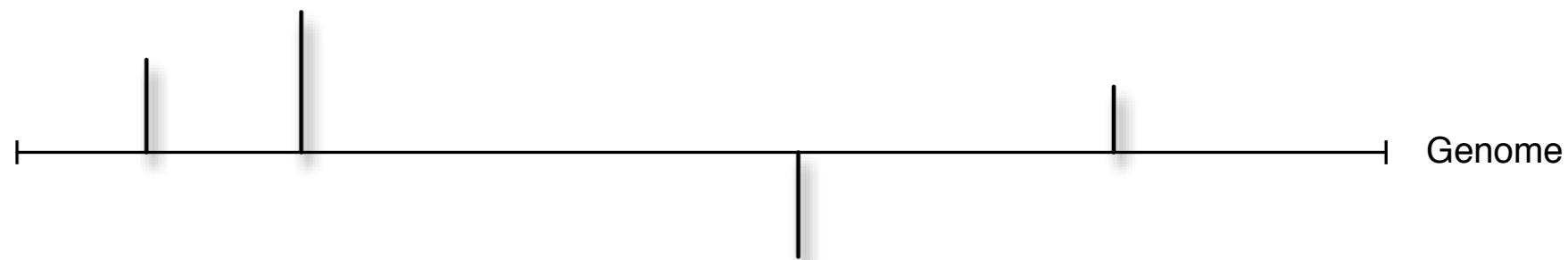


Track
type #4:

Valued
points

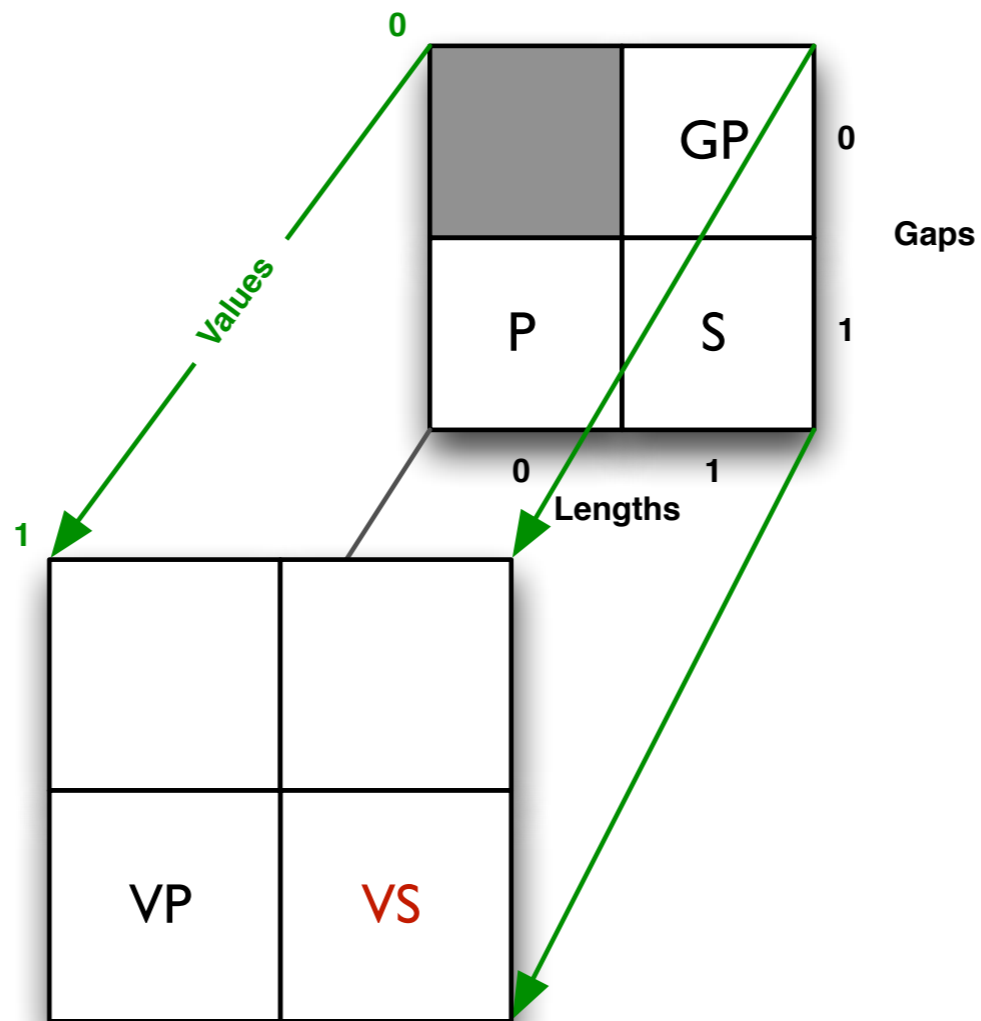


Track
geometry

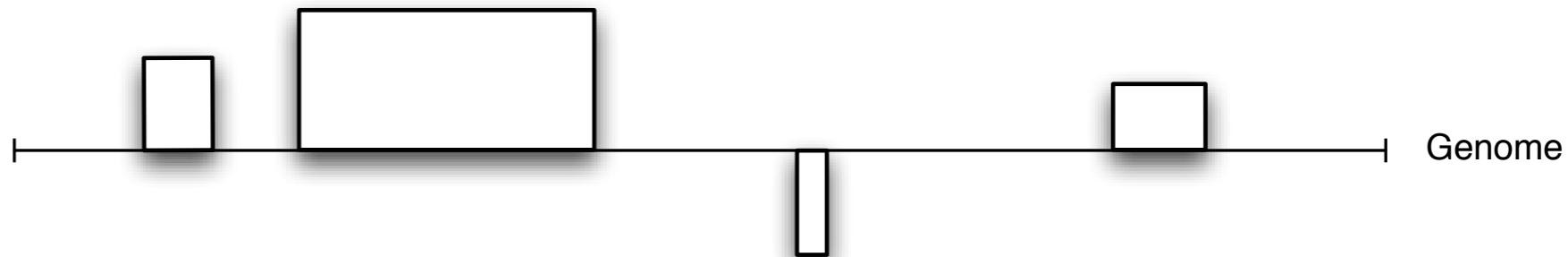


Track
type #5:

Valued
segments

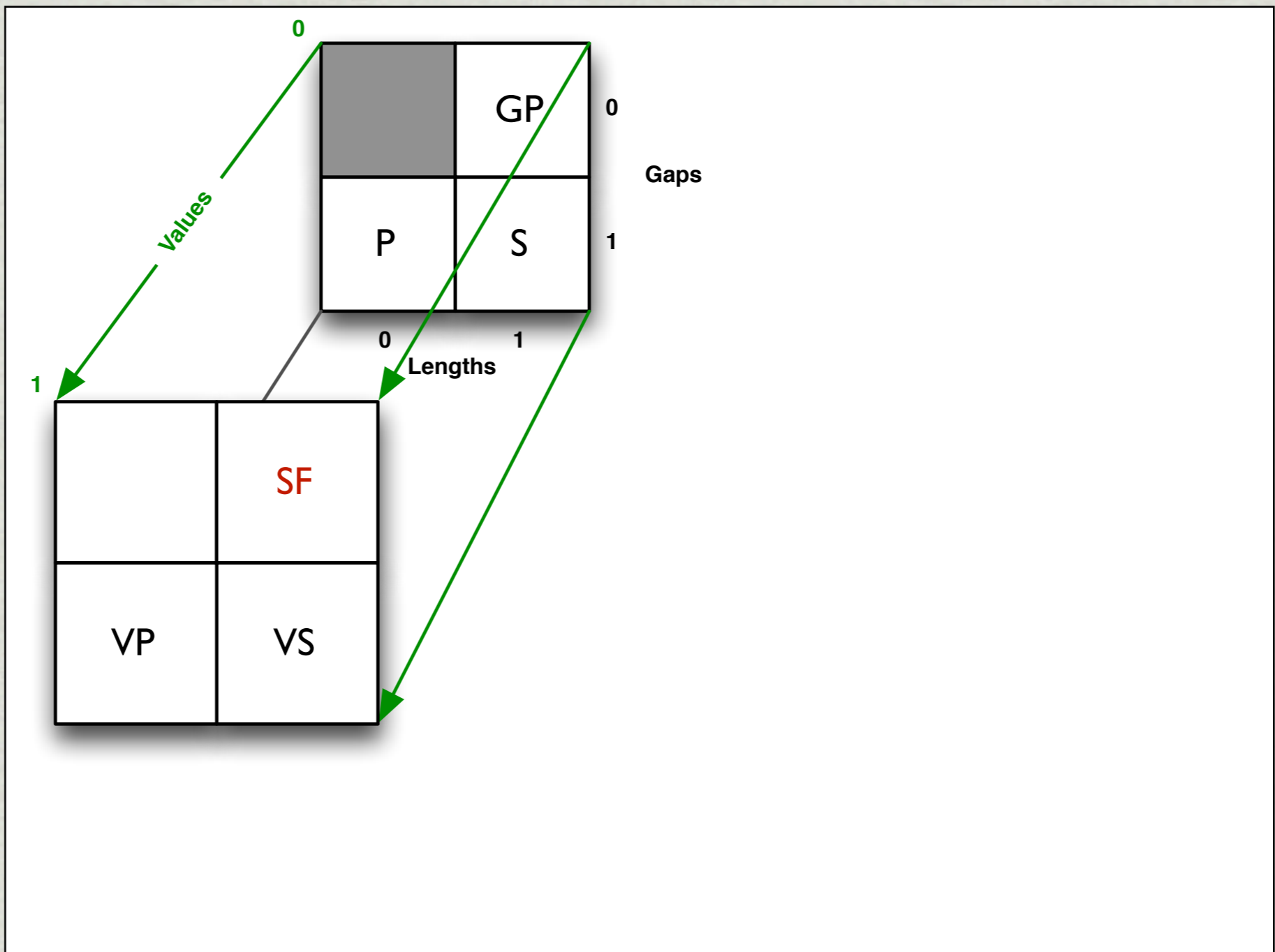


Track
geometry

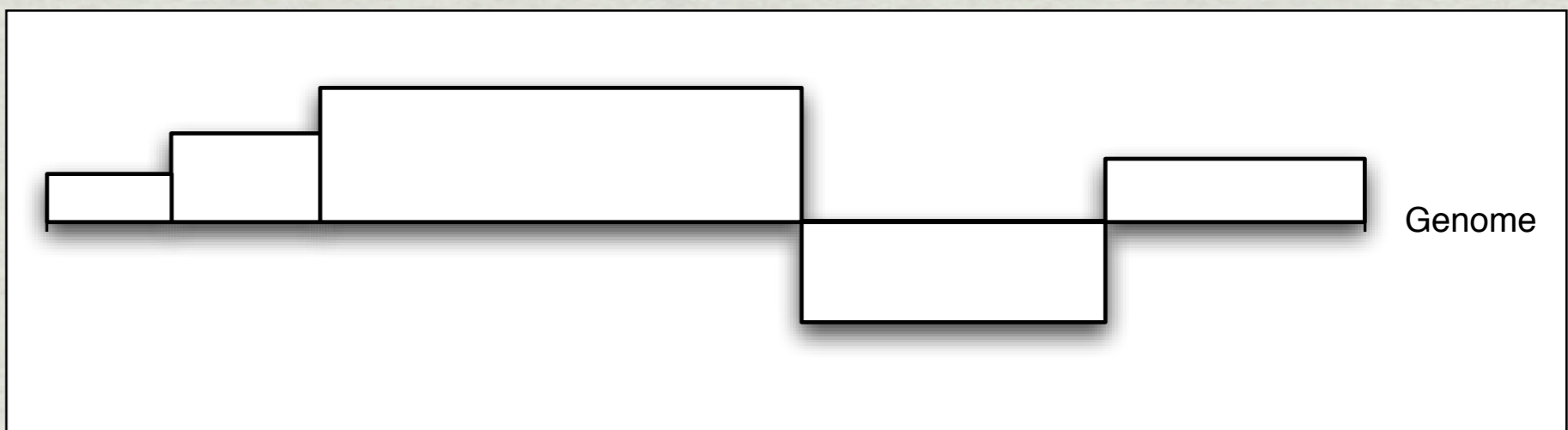


Track
type #6:

Step
function

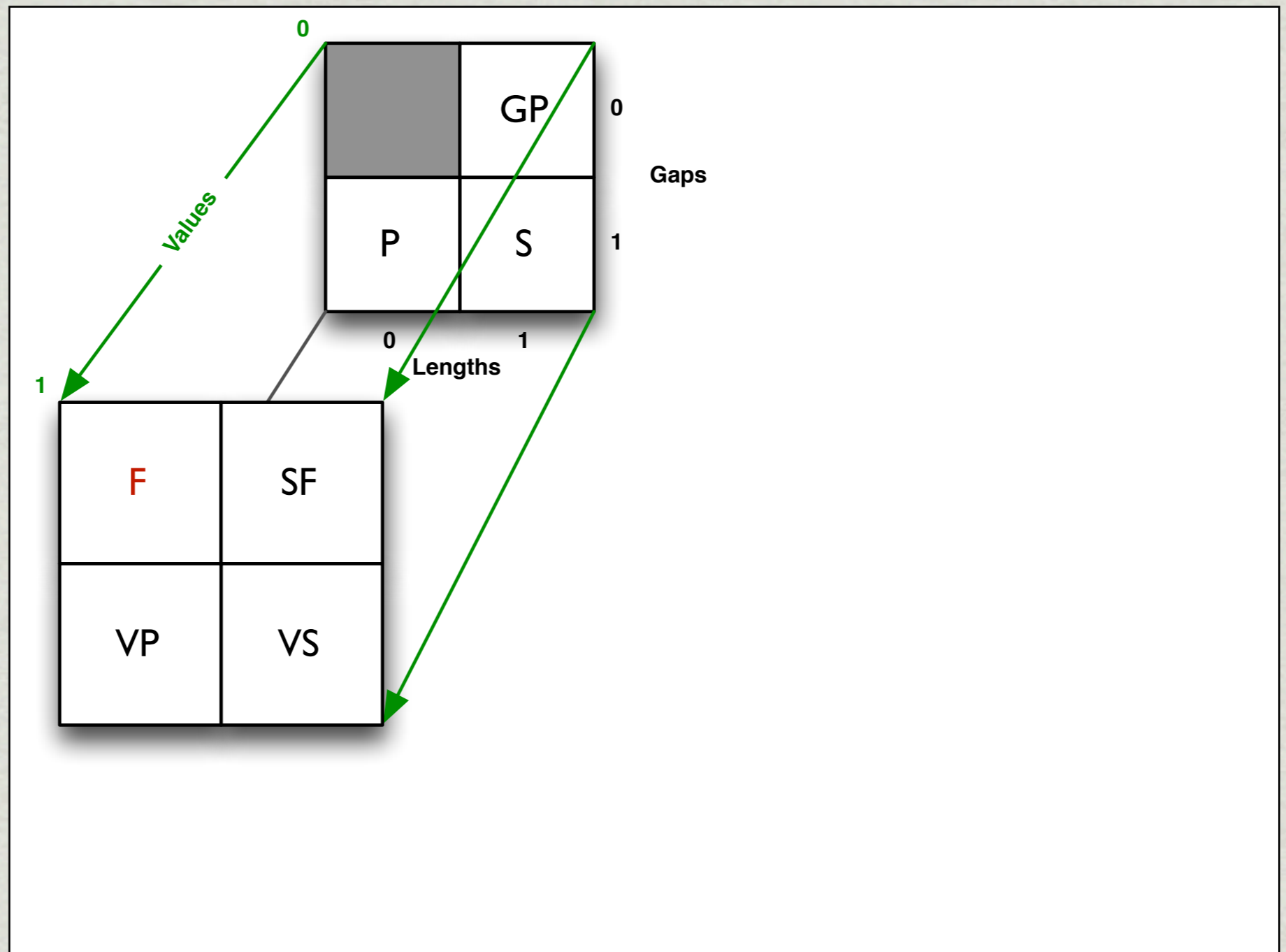


Track
geometry

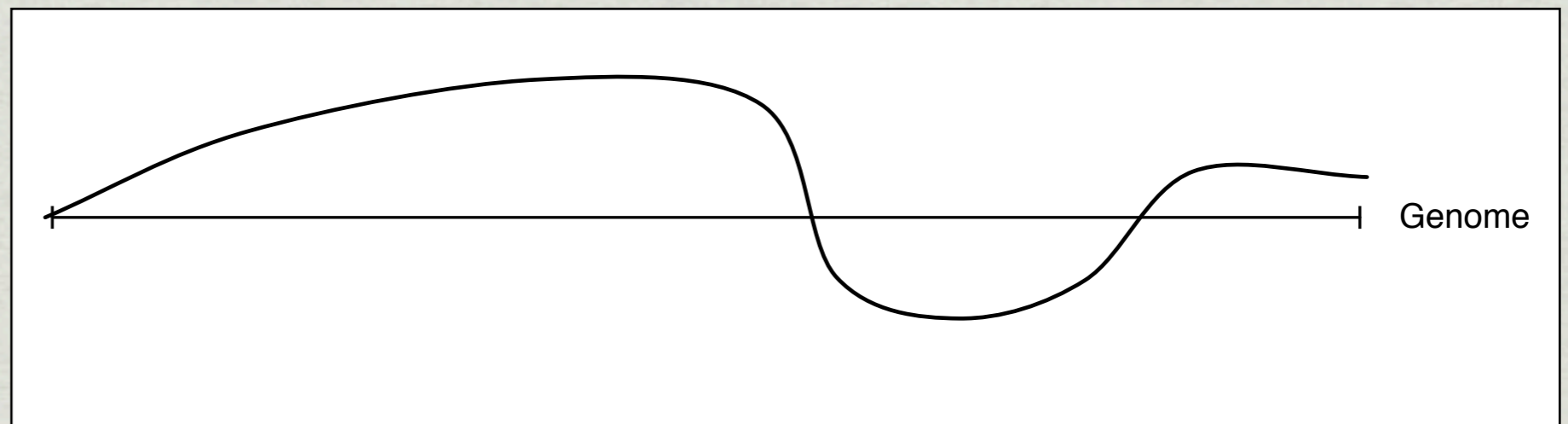


Track
type #7:

Function

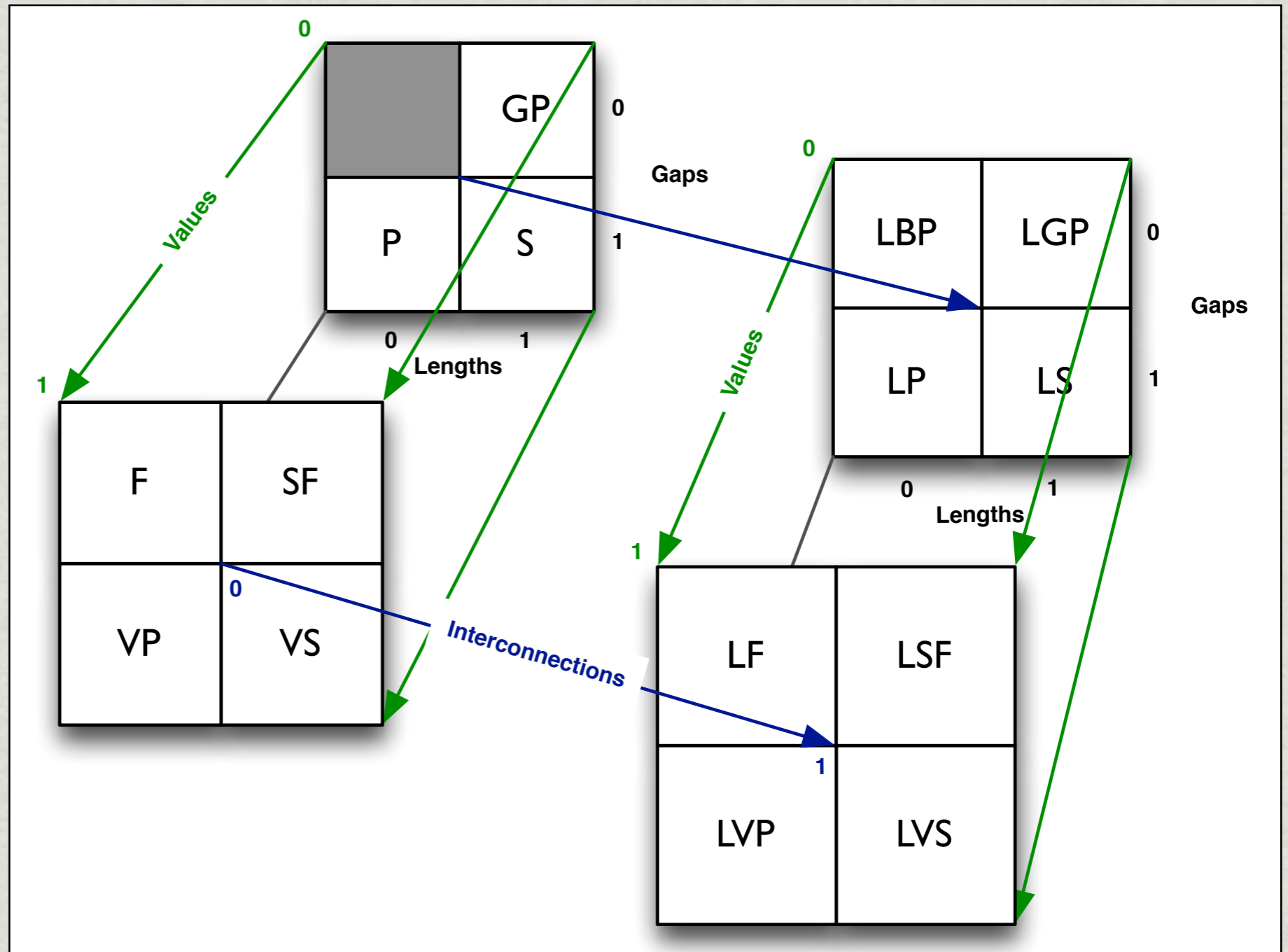


Track
geometry

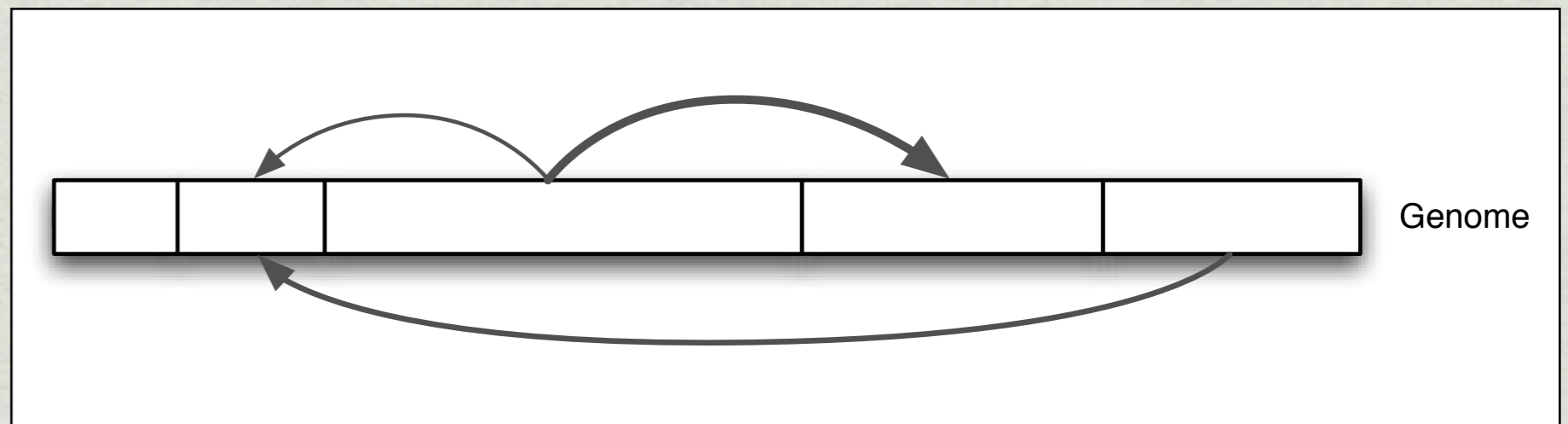


Track
type #8-15:

Linked
track
types

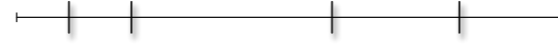


Track
geometry



All 15 track types

Basic track types



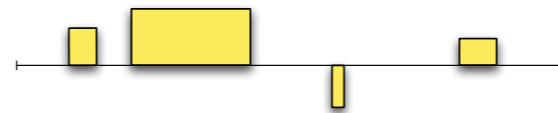
Points (P)



Valued Points (VP)



Segments (S)



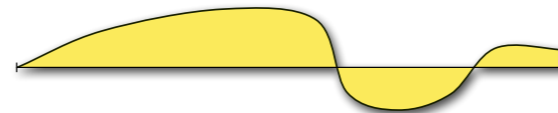
Valued Segments (VS)



Genome Partition (GP)

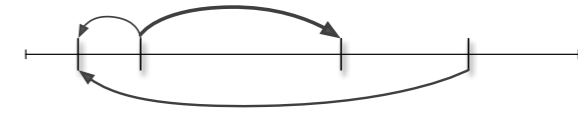


Step Function (SF)

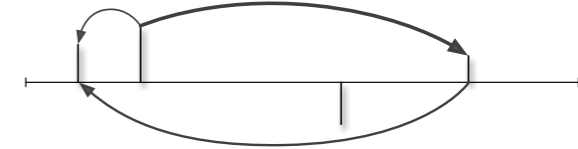


Function (F)

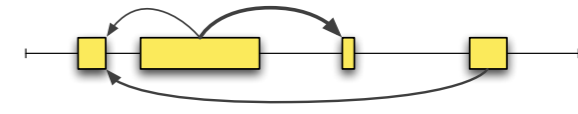
Extended track types



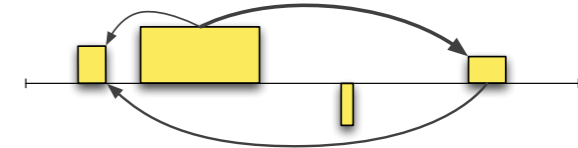
Linked Points (LP)



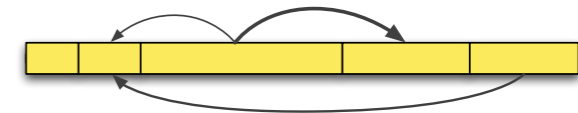
Linked Valued Points (LVP)



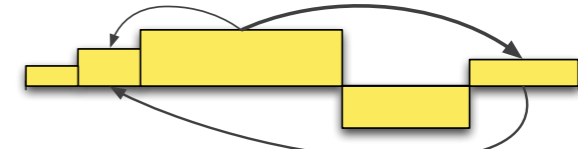
Linked Segments (LS)



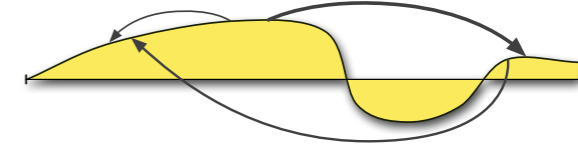
Linked Valued Segments (LVS)



Linked Genome Partition (LGP)



Linked Step Function (LSF)



Linked Function (LF)



Linked Base Pairs (LBP)

Track types relate to analysis questions

P	P	Different frequencies?
P	P	Located nearby?
P	S	Located inside?
P	S	Located <u>nonuniformly</u> inside?
P	S	Located nearby?
S	S	Similar segments?
S	S	Overlap?
S	S	Located nearby?
F	F	Correlated?
P	F	Higher values at locations?
S	F	Higher values inside?
P	VS	Located in segments with high values?
S	VP	Higher values inside segments?
VP	VP	Nearby values similar?
P	VS (c/c)	Located in case segments
VS (c/c)	S	Preferential overlap?
VP (cat)	VS (cat)	Category pairs differentially co-located?
LGP	P	Colocalized in 3D?

Track types supported by different formats

Format	Ref.	Data	Repr.	P	S	VP	VS	GP	SF	F	L	Strand	#Cols	Value type
GFF3/GTF	[2]	General	Tab.	√ ⁽¹⁾	√	√ ⁽¹⁾	√				⁽²⁾	√	9	Float ⁽³⁾
BED/bigBed	[4]	General	Tab./ Bin.	√ ⁽¹⁾	√	√ ⁽¹⁾	√				⁽²⁾	√	3-12	Int(0-1000) /string ⁽⁴⁾
BED15	[4]	Microarray	Tab.			√ ⁽¹⁾	√				⁽²⁾	√	15	List of floats ⁽⁵⁾
bedGraph	[4]	General	Tab.			√ ⁽¹⁾	√						4	Float
WIG/bigWig (fixedStep)	[8]	General	Tab./ Bin.			√	√		√	√			1	Float
WIG/bigWig (variableStep)	[8]	General	Tab./ Bin.			√	√						2	Float
CNT	[36]	Copy number	Tab.			√							4	Float
Personal Genome SNP	[4]	Variation	Tab.			√ ⁽¹⁾	√						7	String ⁽⁶⁾
VCF	[37]	Variation	Tab.			√	√						≥ 8	String ⁽⁶⁾ ⁽³⁾
GVF	[6]	General/ Variation	Tab.	√ ⁽¹⁾	√	√ ⁽¹⁾	√				⁽²⁾	√	9	Float ⁽³⁾
PSL	[4]	Alignment	Tab.		√		√					√	21	Int ⁽⁷⁾
SAM/BAM	[38]	Alignment	Tab./ Bin.		√		√					√	11	Int /string ⁽⁸⁾
BioHDF	[39]	Alignment	Bin.		√		√					√	11	Int /string ⁽⁸⁾
MAF	[4]	Multiple Alignment	Tab.		√		√				⁽⁹⁾	√	2-7	Float /string ⁽⁸⁾
FASTA	[40]	Sequence	Text							√			N/A	Char
DAS XML	[12]	General	XML	√ ⁽¹⁾	√	√ ⁽¹⁾	√				⁽²⁾	√	N/A	Float
BioXSD 1.0	[16]	General	XML	√ ⁽¹⁰⁾	√ ⁽¹⁰⁾	√ ⁽¹⁰⁾	√ ⁽¹⁰⁾				√ ⁽¹¹⁾	√	N/A	Float ⁽¹²⁾
USeq	[19]	General	Bin.	√	√	√	√					√	N/A	Int/float/string
Genomedata	[41]	General	Bin.			√	√		√	√			N/A	Int/float/char

Why all the formats?

1. Different columns
2. Underlying structural differences
- 3. Format-specific characteristics**
4. Tool-specific format requirements
5. Performance issues

Format-specific characteristics

- GFF format:
 - 1-based coordinates
 - Circular genomes
 - Format-specific headers

Why all the formats?

1. Different columns
2. Underlying structural differences
3. Format-specific characteristics
- 4. Tool-specific format requirements**
5. Performance issues

Tool-specific format requirements

- Tools may have specific format requirements

Tool-specific format requirements

- Tools may have specific format requirements
- Galaxy tool (BedTools):
 - intersect multiple sorted BED files

Tool-specific format requirements

- Tools may have specific format requirements
- Galaxy tool (BedTools):
 - intersect multiple sorted BED files
- Segtools toolkit:
 - variant of BED where column 4 is interpreted as a category, rather than a name

Why all the formats?

1. Different columns
2. Underlying structural differences
3. Format-specific characteristics
4. Tool-specific format requirements
5. Performance issues

Performance issues

- Binary (indexed) versions of tabular formats are commonly created to improve performance:
 - SAM => BAM
 - BED => bigBed
 - WIG => bigWig
 - Any interval => Tabix

Introducing GTrack v1.0...

GTrack example (simple)

chr1	2396586	2827369
chr1	92014277	93306499
chr1	100983315	101455310
chr1	116832232	116909542
chr1	190733439	190814781
chr1	199128354	199336605
...		

GTrack example (intermediate)

```
##gtrack version: 1.0
##track type: valued segments
##value type: binary
##sorted elements: true
##1-indexed: true
##end inclusive: true
##my custom header: Yes!
###seqid      start      end      value
chr1         2396587   2827369   1
chr1         92014278  93306499  1
chr1        100983316 101455310  0
chr1        116832233 116909542  0
chr1        190733440 190814781  1
chr1        199128355 199336605  0
...
```

GTrack example (intermediate)

```
##gtrack version: 1.0
##track type: valued segments
##value type: binary
##sorted elements: true
##1-indexed: true
##end inclusive: true
##my custom header: Yes! Flexible columns
###seqid      start      end      value
chr1      2396587    2827369    1
chr1      92014278   93306499   1
chr1      100983316  101455310  0
chr1      116832233  116909542  0
chr1      190733440  190814781  1
chr1      199128355  199336605  0
...
```

GTrack example (intermediate)

```
##gtrack version: 1.0
```

```
##track type: valued segments
```

```
##value type: binary
```

```
##sorted elements: true
```

```
##1-indexed: true
```

```
##end inclusive: true
```

```
##my custom header: Yes!
```

```
###seqid      start      end      value
```

```
chr1          2396587    2827369    1
```

```
chr1          92014278   93306499   1
```

```
chr1         100983316  101455310  0
```

```
chr1         116832233  116909542  0
```

```
chr1         190733440  190814781  1
```

```
chr1         199128355  199336605  0
```

```
...
```

Supports all 15
track types

GTrack example (intermediate)

```
##gtrack version: 1.0
##track type: valued segments
##value type: binary
##sorted elements: true
##1-indexed: true
##end inclusive: true
##my custom header: Yes!
```

Format-specific characteristics, custom headers

##seqid	start	end	value
chr1	2396587	2827369	1
chr1	92014278	93306499	1
chr1	100983316	101455310	0
chr1	116832233	116909542	0
chr1	190733440	190814781	1
chr1	199128355	199336605	0
...			

GTrack example (intermediate)

```
##gtrack version: 1.0
##track type: valued segments
##value type: binary
##sorted elements: true
##1-indexed: true
##end inclusive: true
##my custom header: Yes!
###seqid      start      end      value
chr1          2396587   2827369   1
chr1          92014278  93306499  1
chr1          100983316 101455310 0
chr1          116832233 116909542 0
chr1          190733440 190814781 1
chr1          199128355 199336605 0
...
```

Tool-specific format requirements

GTrack example (intermediate)

Specification of GTrack subtypes via URL with automatic validation

```
##subtype url: http://my.edu/frmt.gtrack  
chr1      2396587      2827369      1  
chr1      92014278     93306499     1  
chr1      100983316    101455310    0  
chr1      116832233    116909542    0  
chr1      190733440    190814781    1  
chr1      199128355    199336605    0  
...
```

GTrack example (advanced)

```
#Example track
##gtrack version: 1.0
##track type: linked genome partition
##edge weights: true
###end      id      edges
####seqid=chr1;start=1000;end=2000
1015      A      C=1.3,D=0.1
1060      B      A=1.0
1154      C      A=1.3
1267      D      .

...
```

GTrack example (advanced)

```
#Example track
##gtrack version: 1.0
##track type: linked genome partition
##edge weights: true
###end      id      edges
####seqid=chr1;start=1000;end=2000
1015      A      C=1.3,D=0.1
1060      B      A=1.0
1154      C      A=1.3
1267      D      .
...

```

Specification of
edges with weights

GTrack example (advanced)

```
#Example track
##gtrack version: 1.0
##track type: linked genome partition
##edge weights: true
###end      id      edges
####seqid=chr1;start=1000;end=2000
1015      A      C=1.3,D=0.1
1060      B      A=1.0
1154      C      A=1.3
1267      D      .
...
```

**Specification of
bounding regions
(to exclude missing
data / assembly gaps)**

GTrack performance

- GTrackCore library is in development, parsing GTrack files into a binary format for fast analysis

GTrack performance

- GTrackCore library is in development, parsing GTrack files into a binary format for fast analysis
- Based on the backbone of the Genomic HyperBrowser

GTrack performance

- GTrackCore library is in development, parsing GTrack files into a binary format for fast analysis
- Based on the backbone of the Genomic HyperBrowser
- Blazingly fast, vector-based calculations using numpy!

GTrack performance

- GTrackCore library is in development, parsing GTrack files into a binary format for fast analysis
- Based on the backbone of the Genomic HyperBrowser
- Blazingly fast, vector-based calculations using numpy!
- Converts to/from many common formats

GTrack performance

- GTrackCore library is in development, parsing GTrack files into a binary format for fast analysis
- Based on the backbone of the Genomic HyperBrowser
- Blazingly fast, vector-based calculations using numpy!
- Converts to/from many common formats
- Planned to move to PyTables:
 - HDF5-based, compressed storage
 - advanced indexing and querying
 - solid performance
 - several tracks in one file

Galaxy relevance of GTrack

- 15 track types (including 3D type data)!

Galaxy relevance of GTrack

- 15 track types (including 3D type data)!
- Provides easy way to handle tool-specific format requirements

Galaxy relevance of GTrack

- 15 track types (including 3D type data)!
- Provides easy way to handle tool-specific format requirements
- Everything specified in the file itself, interoperable with other systems

Galaxy relevance of GTrack

- 15 track types (including 3D type data)!
- Provides easy way to handle tool-specific format requirements
- Everything specified in the file itself, interoperable with other systems
- GTrack fully integrated with the HyperBrowser, along with a set of GTrack-specific tools

Galaxy relevance of GTrackCore

- HDF5-based binary file as history element

Galaxy relevance of GTrackCore

- HDF5-based binary file as history element
- Fast performance (after indexing)

Galaxy relevance of GTrackCore

- HDF5-based binary file as history element
- Fast performance (after indexing)
- Library of operations, analysis, statistics to ease tool development

Interested?

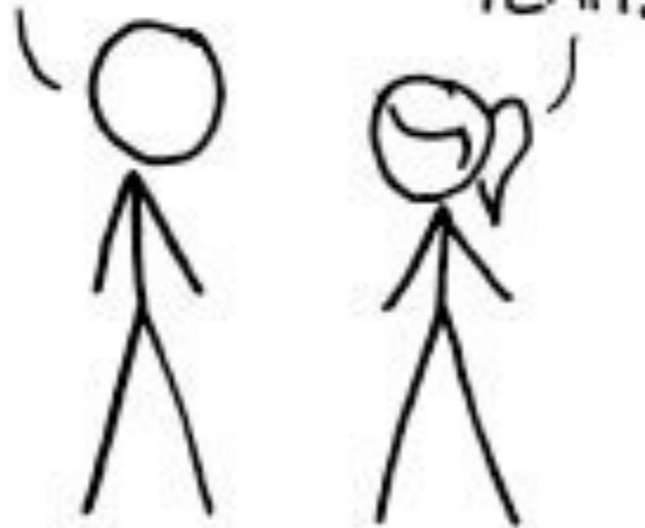
We need more developers!

HOW STANDARDS PROLIFERATE:

(SEE: A/C CHARGERS, CHARACTER ENCODINGS, INSTANT MESSAGING, ETC.)

SITUATION:
THERE ARE
14 COMPETING
STANDARDS.

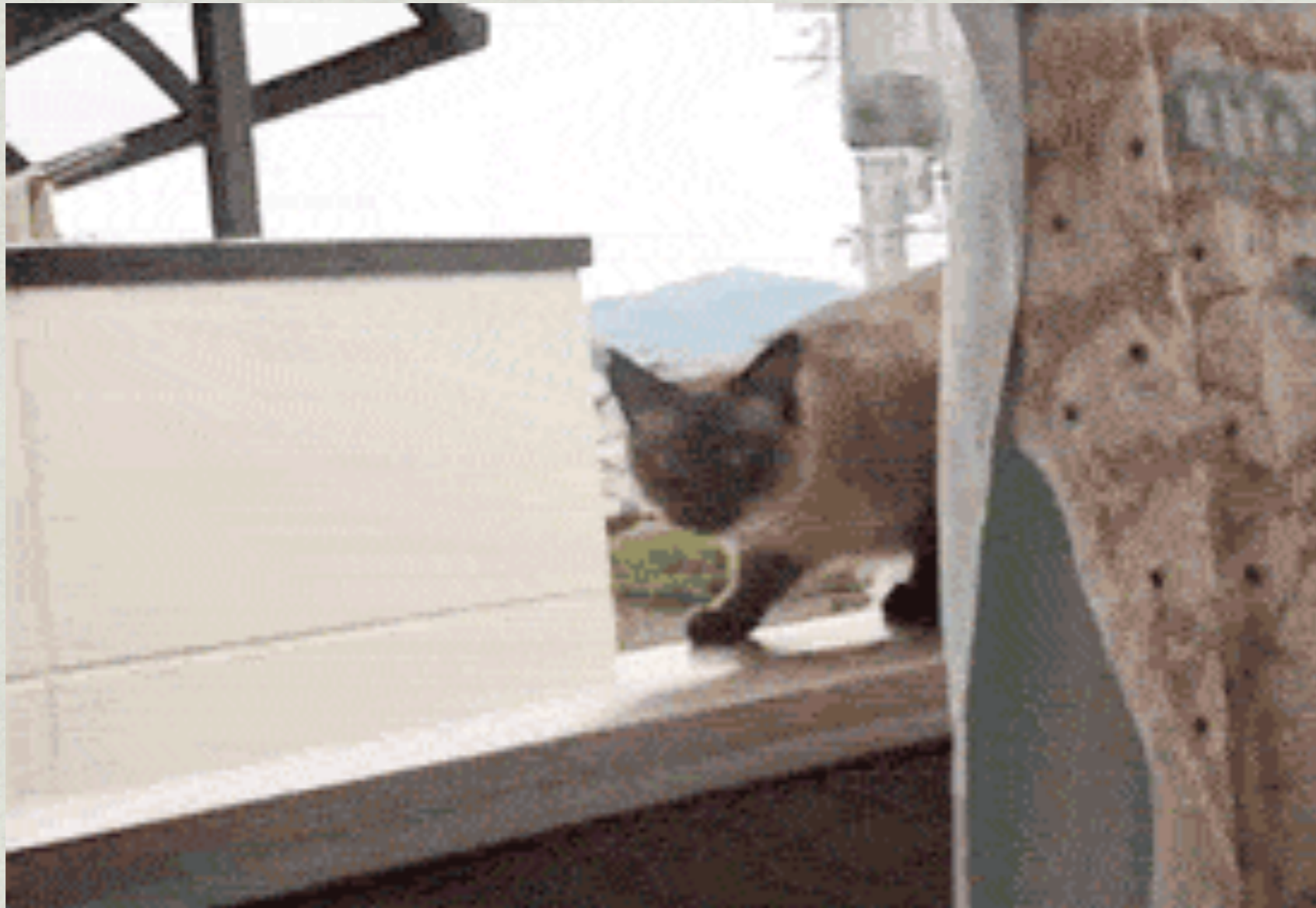
14?! RIDICULOUS!
WE NEED TO DEVELOP
ONE UNIVERSAL STANDARD
THAT COVERS EVERYONE'S
USE CASES.



SOON:

SITUATION:
THERE ARE
15 COMPETING
STANDARDS.

**Nothing ventured,
nothing gained...**



Acknowledgements

Sveinung Gundersen

Matúš Kalaš

Osman Abul

Marcin Cieřlik

Arnoldo Frigessi

Eivind Hovig

Geir Kjetil Sandve

More info

Gundersen S, Kalas M, Abul O, Frigessi A, Hovig E, Sandve GK:
Identifying elemental genomic track types and representing them uniformly. BMC Bioinformatics 2011, 12:494.

Specification available at:
www.gtrack.no

GTrackCore source code:
<https://github.com/i000/gtrackcore>