

Sequence Analysis of Human Immunodeficiency Virus Type 1

Stephanie Lucas ^{1,2}

Mentor: Panayiotis V. Benos ^{1,3}

With help from: David L. Corcoran ⁴

¹ Bioengineering and Bioinformatics Summer Institute, Department of Computational Biology, University of Pittsburgh

² Department of Biology, University of San Francisco

³ Department of Computational Biology, University of Pittsburgh

⁴ Departments of Computational Biology and Human Genetics, University of Pittsburgh

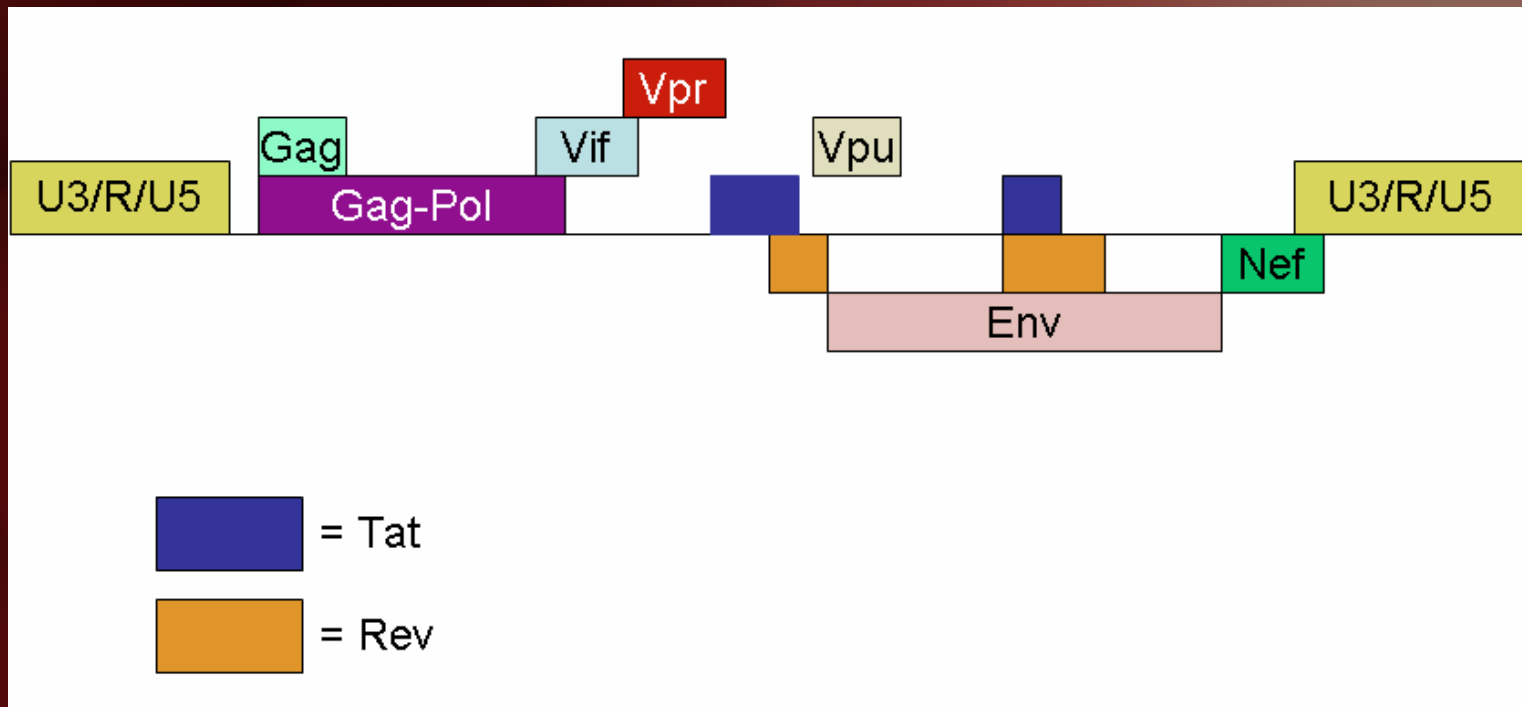
- Some Facts as of 2005...
 - There are 38.6 M people in the world living with HIV/AIDS
 - US has 1.2 M out of the 1.3 M living in N. America.
 - Kenya has 1.3 M, and sub-saharan Africa 24.5 M
- There is currently no successful therapy.
Production of vaccine is difficult due to the high mutation rate of the virus.
The effects are severe both economically (cost of care) and socially .

Purpose

- We will study the evolution of different parts of the HIV-1 genome
- Parts that evolve slower could act as potential vaccine targets

HIV-1 Genome

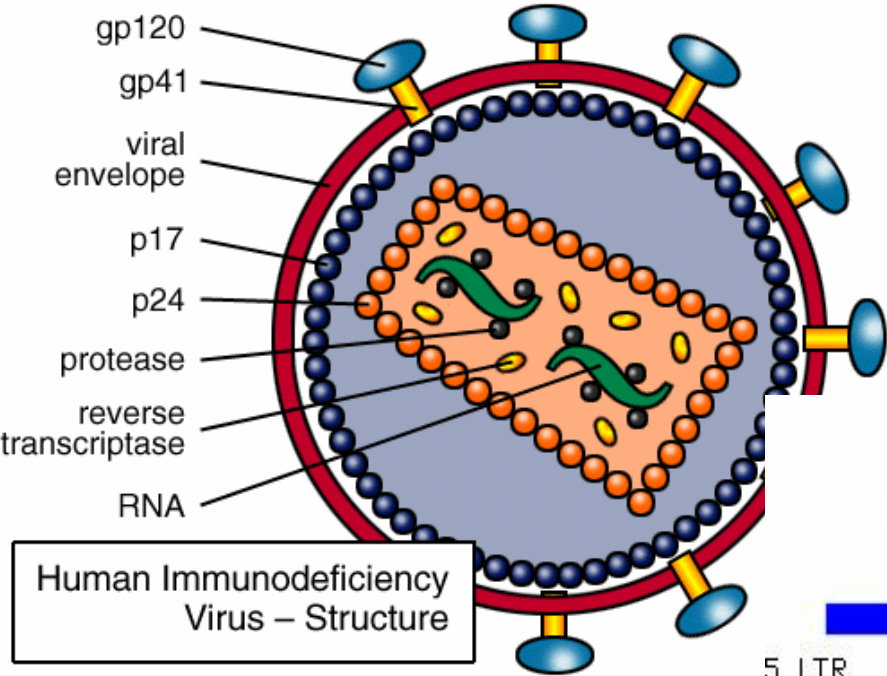
- Reference Sequence is 9181 base pairs long
- Contains 9 genes



- Gag: codes for internal structural proteins and capsid proteins
- Gag-Pol: codes for the three enzymes necessary for replication Vpr
- Vpu: virus protein U
- Tat: transactivator protein
- Rev: regulator of expression of virus protein

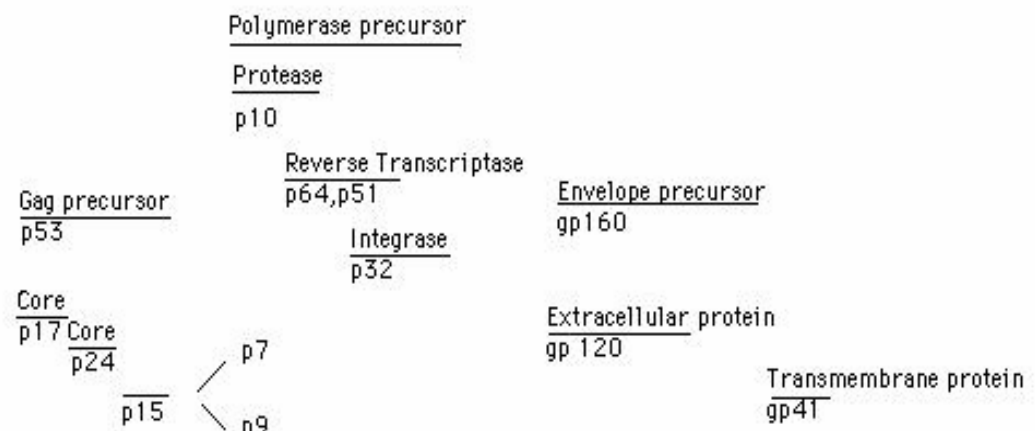
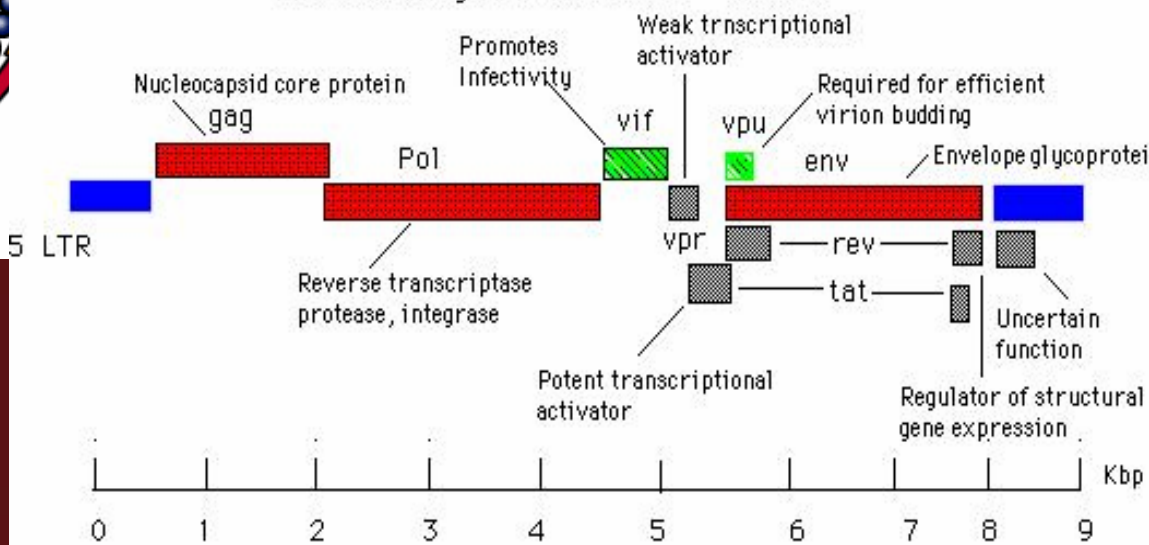
The genes cont'd

- Env: codes for the surface proteins gp120 and gp41 that protrude from the lipid envelope and attach to cellular receptors
- Nef: an enhancing factor
- Vpr: virus protein R
- Vif: virus infectivity factor

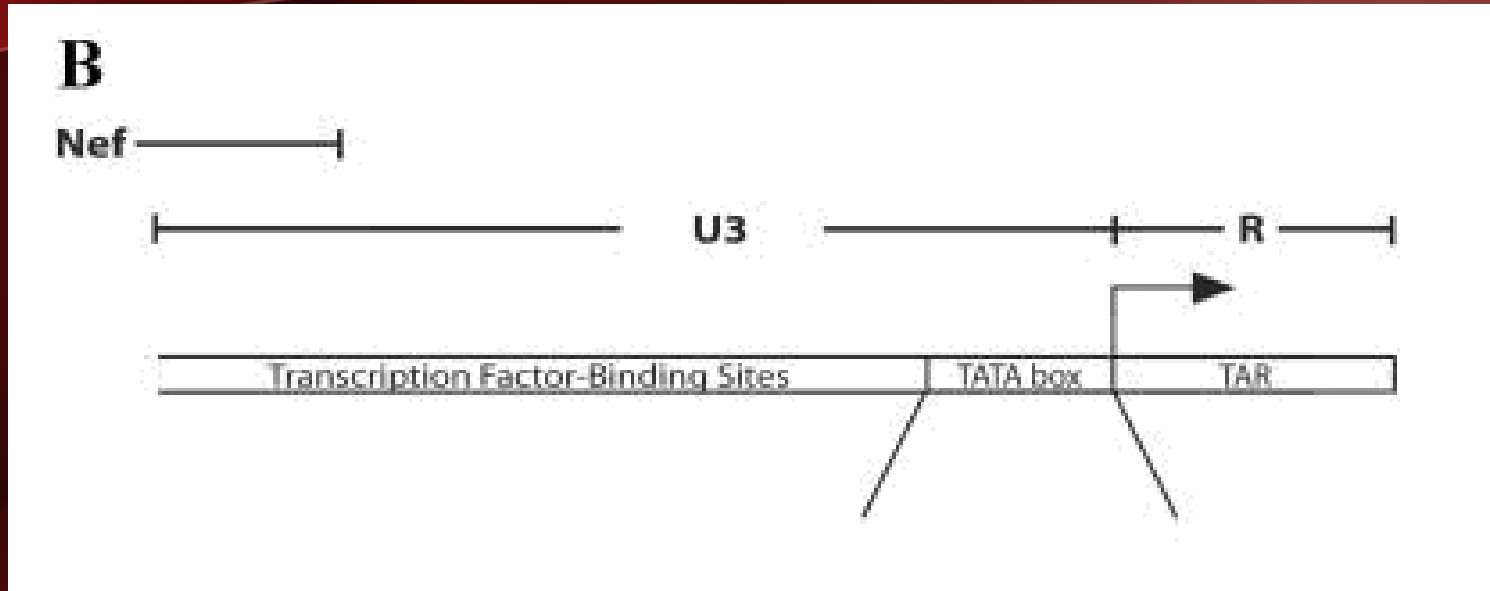


Human Immunodeficiency Virus – Structure

Genetic Organisation of HIV-1 virus

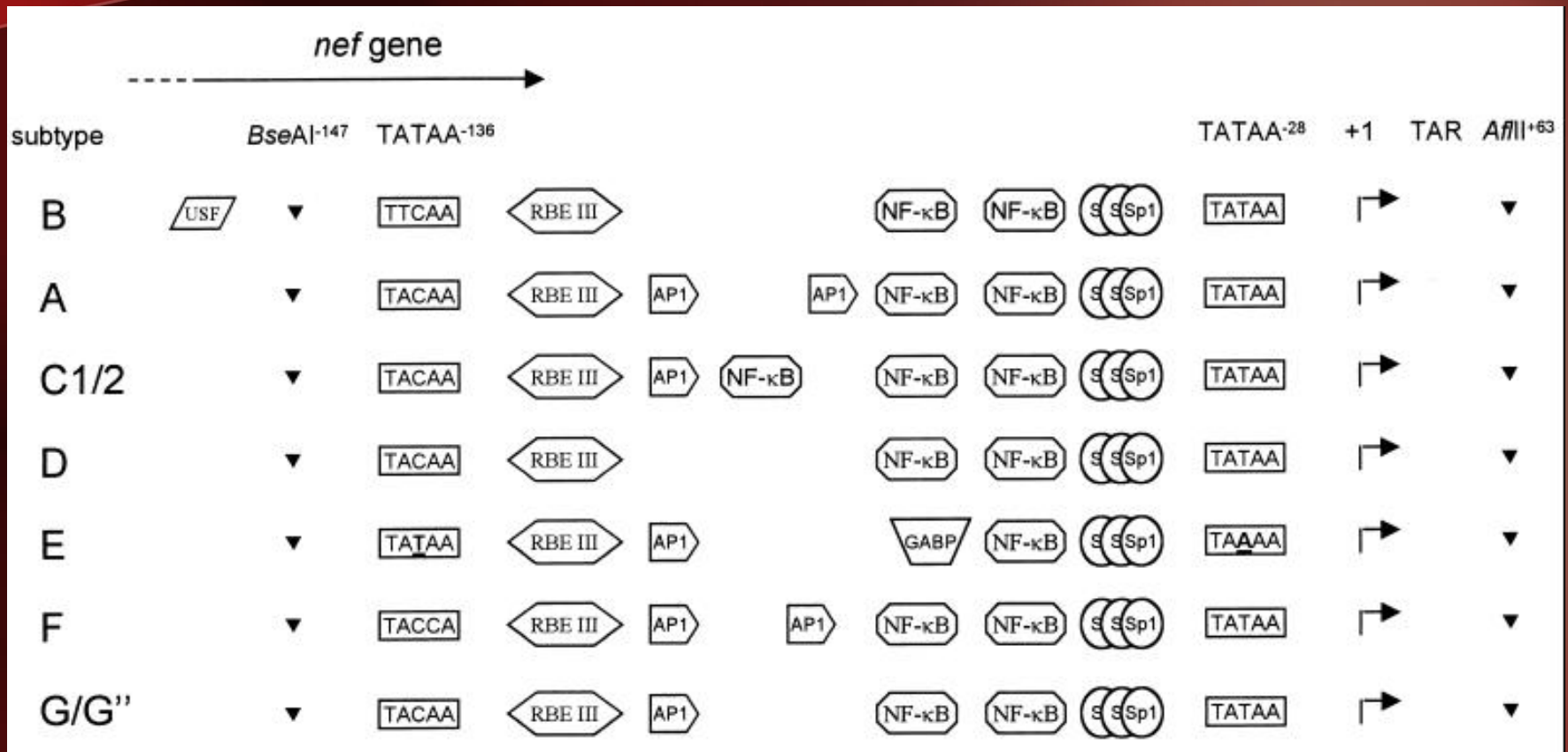


U3_R_U5 Region



- The U3_R region is the regulatory region and contains most of the transcription factor binding sites

Predicted TFBS



Overview of Research---Our Goals

- To track if there are any differential selective pressure on parts of the genome
- Identify regions of higher/lower variability
- Predict and confirm TFBS within the promoter region

Obtaining the sequences

- a. Looked up the Ref Seq from the database
- b. Searching in the public databases yielded 1,183 genomes
- c. Split the Ref Seq into individual genes and regulatory regions
 - coding/ regulatory regions only
- d. Removed overlapping sequences and Start/Stop codons
 - There are differential constrains within individual bases
 - As a consequence, 2 genes were not analyzed- TAT and GAG
 - Start/ Stop codons are relatively invariable and may stray the conserved sequence count
- e. Did a BLAST search against the 1,183 genomes to extract out each gene from the sequences and remove identical sequences- left with about 200 sequences
- f. Align with Clustal W using the MEGA software package

The sequences were then ready for analysis...

Example CLUSTAL W

Alignment Explorer (C:\Documents and Settings\lucas\Desktop\SET 1- aligned MEGA sequences- ALL sequences, even prolematic ones\align 2- MASWIF-alignment 2.mas)

Data Edit Search Alignment Web Sequencer Display Help



DNA Sequences | Translated Protein Sequences

gi 4205033 gb U6958	GGATGCTAAA	---	TTGGTAA	TAA	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTGGG	CCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 60651856 gb AY83	GGATGCTAAA	---	TTGGTAA	TAA	TGCAAC	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTGGG	CCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 60651905 gb AY83	GGATGCTAGA	---	TTGGTAA	TAA	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTGGG	CCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 60652082 gb AY83	GGATGCTAGA	---	TTGGTAA	TAA	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTGGG	CCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 60652044 gb AY83	GGATGCTAAA	---	TTGGTAG	TAA	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTAGG	TCAGGG	AGTCTCC	---	GTAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 10436130 gb AF25	GGATGCAAAA	---	TTGGTAA	TAA	CAACA	TATTTGGGG	TCTA	CATACA	GGAGAAA	GAGACTGG	CATTTGGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 3511259 gb AF075	GGATGCTAGA	---	TTGGTAA	TAA	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTAGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 55275205 gb AY77	GGATGCTAGA	---	TTGGTAG	TAA	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTGGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 60651886 gb AY83	GGATGCTAGA	---	TTGGTAA	TAA	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTGGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 60652112 gb AY83	GGATGCTAGA	---	TTGGTAG	TAA	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTGGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 3193272 gb AF069	GGATGCTAGC	---	TTGGTAG	TAA	CAACA	TATTTGGGG	TCTA	CATACA	GGAGAAA	GAGACTGG	CATTTGGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 60544776 gb AY83	GGATGCTAGA	---	TTGGTAG	TAA	CAACA	TATTTGGGG	TCTGC	CACAC	GGAGAAA	GAGACTGG	CATTTGGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	CAAA	AGATATA	GCACAC	AAAGTAGA	
gi 10436120 gb AF25	GGATGCAAAA	---	TTGGTAA	TAA	CAACA	TATTTGGGG	TCTA	CATACA	GGAGAAA	GACTGG	CATTTGGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 328030 gb M17449	GGATGCTAGA	---	TTGGTAA	TAA	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTAGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 4205051 gb U6959	GGATGCTAAA	---	TTGGTAA	TAA	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTGGG	CCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 55275225 gb AY77	GGATGCTAGA	---	TTGGTAG	TAA	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTGGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 60651827 gb AY83	GGATGCTAGA	---	TTGGTAA	TAA	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTGGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 60544786 gb AY83	GGATGCTAGA	---	TTGGTAA	TAA	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTGGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 328565 gb M17451	GGATGAAA	GG	---	TTGGTAA	TAA	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTGGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA
gi 3098582 gb AF049	GGATGCTAGA	---	TTGGTAG	TAA	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTGGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 29119285 gb AY17	GGAAAGCTAGA	---	TTAGTAG	TAA	CAACA	TATTTGGGG	TCTGC	AAAC	GGAGAAA	GAGACTGG	CATTTGGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 4204988 gb U6958	GGATGCTAAA	---	TTGGTAA	TAG	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTGGG	CCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 55925120 gb AY81	GGATGCTAAA	---	TTGGTAA	TAG	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTGGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 37677783 gb AY33	GGATGCTAGA	---	TTGGTAG	TAA	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTGGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 37677793 gb AY33	GGATGCTAGA	---	TTGGTAG	TAA	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTGGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 60652014 gb AY83	GGATGCTAGA	---	TTAGTAA	TAA	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTGGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 190960709 dbj AB2	GGATGCGAAA	---	TTGGTAA	TAA	CAACA	TATTTGGGG	TCTA	CATACA	GGAGAAA	GAGACTGG	CATTTGGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 10436111 gb AF25	GGATGCGAAA	---	TTGGTAA	TAA	CAACA	TATTTGGGG	TCTA	CATACA	GGAGAAA	GAGACTGG	CATTTGGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 29119265 gb AY17	GGATGCTAGA	---	TTGGTAA	TAA	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTGGG	TCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 4205042 gb U6959	AGATGCTAAAT	---	TTGGTAA	TAA	CAACA	TATTTGGGG	TCTGC	ATTCA	GGAGAAA	GAGACTGG	CATTTGGG	CCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 4205069 gb U6959	AGATGCTAAAT	---	TTGGTAA	TAA	CAACA	TATTTGGGG	TCTGC	ATTCA	GGAGAAA	GAGACTGG	CATTTGGG	CCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	
gi 55925112 gb AY81	GGATGCTAAA	---	TTGGTAA	TAG	CAACA	TATTTGGGG	TCTGC	TATACA	GGAGAAA	GAGACTGG	CATTTGGG	CCAGGG	AGTCTCC	---	ATAGAA	TGG	AGG	AAAA	AGATATA	GCACAC	AAAGTAGA	

Site # 156 with w/o Gaps

Infinite Sites Neutral Model

- Developed by Kreitman and Hudson (1991)
- Focused on neutral (silent) mutations
 - Removes the functional constraints in order to focus on the genetic drift alone
- θ = level of polymorphism $\theta = 4N_e\mu$
 - Actual θ value cannot be found because N and μ are difficult to obtain...

First we had to calculate the number of segregating sites per silent nucleotide (p_s)...

$$p_s = S / n$$

$$\hat{\theta} = p_s / a_1$$

$$\text{Var}(\hat{\theta}) = \text{Var}(p_s) / a_1^2$$

Where....

S = # silent segregating sites

n = # possible silent sites

$$a_k = \sum_{x=1}^{m-1} x^{-k}$$

Calculate the variables...

Example sequence:

	<u>Leu</u>	<u>Gly</u>
Seq #1	CTG	GGC
Seq #2	CTG	GGC
Seq #3	CTA	GGC
Seq #4	CTT	GGC

•Kreitman and Hudson method

- $s = 1$
- $n = 2$
- Counts whether or not a column of the sequence is a silent site

- The 3rd codon position is a potential silent site

Problems...

- Upon further investigation, we discovered that some of the gene sequences were nonfunctional due to mutation or premature stop codons
- We thought that this may influence the results and decided to create 3 sets of data and run each set through the script

Run the test with 3 different sets of data...

SET 1-
All sequences

SET 2-
Problematic
sequences
deleted from
each gene

SET 3-
Problematic
sequences
deleted from all
genes

VPR
sequence

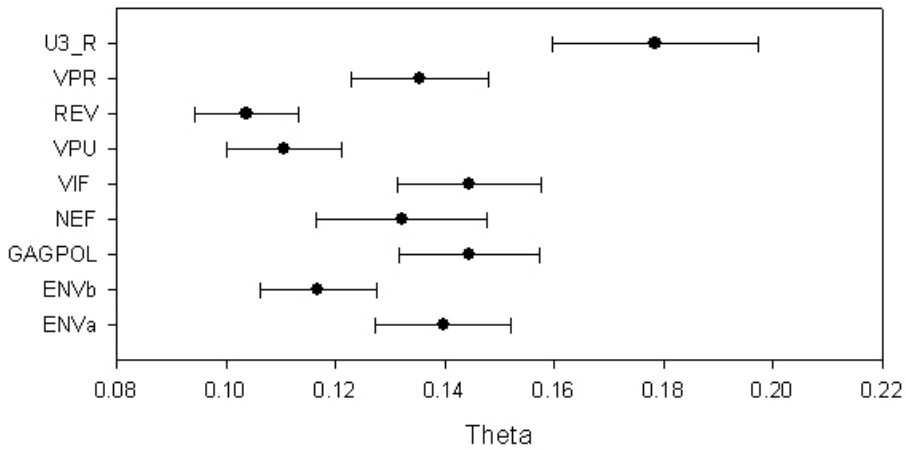
```
graph TD; S1[SET 1- All sequences]; S2[SET 2- Problematic sequences deleted from each gene]; S3[SET 3- Problematic sequences deleted from all genes]; S2 --> A[Delete all short sequences]; S3 --> B[Truncate all sequences];
```

Delete all short
sequences

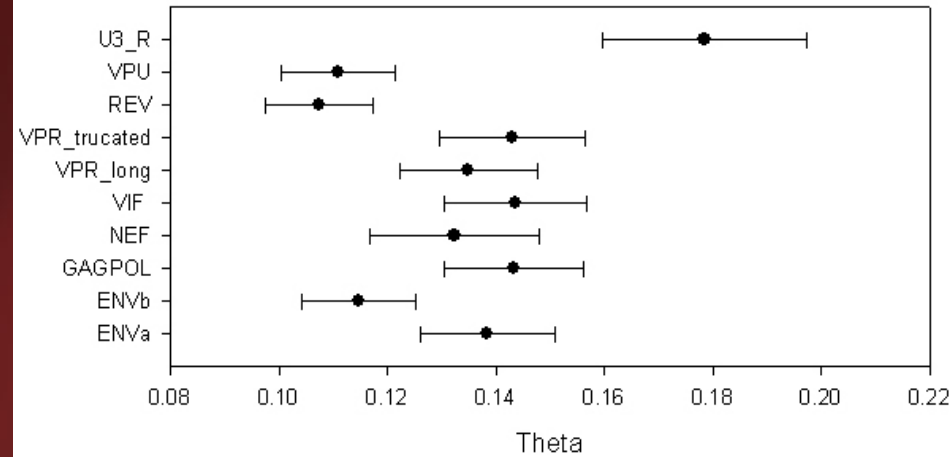
Truncate all
sequences

Results

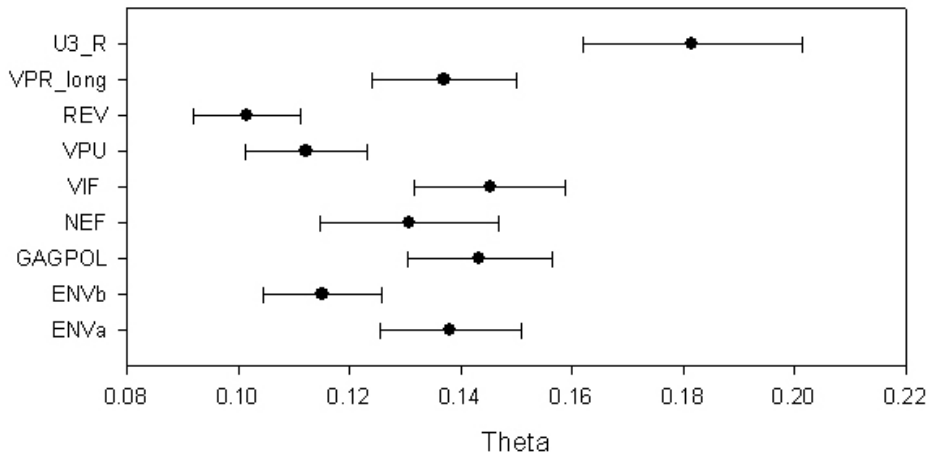
Set 1 of New



Set 2 of New



Set 3 of New



Results Interpreted

- U3_R region is less constrained than the other regions (due to higher Θ values)
- VPU, REV, ENVb are more constrained than other genes (due to lower Θ values)

Comparison with other closely related species

- McDonald and Kreitman (1991)
- The idea is that the ratio of nonsynonymous to synonymous mutations within a species (polymorphisms) should be the same as the ratio between species if the mutations are neutral
- Used HIV-2, SIV-1, and SIV-2 in order to test this
- Although they were close relatives, the sequences were too different, and could not be compared

Instead we used the different subtypes of HIV-1 compared to that of Subtype B

- There are 9 subtypes: A, B, C, D, F, G, H, J, K

Methods:

- Obtained complete genomes for each subtypes
- Did similar extraction methods used in HIV-1 sequence
- Do to limited time, analyzed all sequences of each gene of subtype B and a random sequence of a gene per subtype

Results

	Observed			Expected			
ENV	Fixed	Polymorph	<i>residuals</i>	Fixed	Polymorphi	G-value	
Nonsynonymyn	439	790	1229	434.50	794.50	4.523008465	-4.48702125
Synonymou:	395	735	1130	399.50	730.50	-4.47435055	4.513618736
	834	1525	2359			G-test value	0.150510788
						p-value	0.698047708
GAGPOL	Fixed	Polymorphic					
Nonsynonymyn	219	601	820	266.24	553.76	-42.7794747	49.20393745
Synonymou:	482	857	1339	434.76	904.24	49.72230631	-45.9873991
	701	1458	2159			G-test value	20.3187399
						p-value	6.5555E-06
NEF	Fixed	Polymorphic					
Nonsynonymyn	51	128	179	55.65	123.35	-4.45066615	4.737231073
Synonymou:	46	87	133	41.35	91.65	4.902885968	-4.53059894
	97	215	312			G-test value	1.317703895
						p-value	0.251004736
REV	Fixed	Polymorphic					
Nonsynonymyn	51	66	117	46.58	70.42	4.619723277	-4.27516555
Synonymou:	35	64	99	39.42	59.58	-4.1594237	4.576460189
	86	130	216			G-test value	1.523188436
						p-value	0.21713774
VIF	Fixed	Polymorphic					
Nonsynonymyn	82	134	216	76.29	139.71	5.914699402	-5.58826033
Synonymou:	72	148	220	77.71	142.29	-5.49156862	5.819345174
	154	282	436			G-test value	1.308431244
						p-value	0.252679029
VPR	Fixed	Polymorphic					
Nonsynonymyn	25	67	92	32.34	59.66	-6.4331233	7.770285908
Synonymou:	46	64	110	38.66	71.34	7.992447149	-6.94572213
	71	131	202			G-test value	4.767775269
						p-value	0.02899728
VPU	Fixed	Polymorphic					
Nonsynonymyn	39	53	92	37.88	54.12	1.133974299	-1.10602586
Synonymou:	24	37	61	25.12	35.88	-1.09240413	1.134875083
	63	90	153			G-test value	0.140838782
						p-value	0.707448569

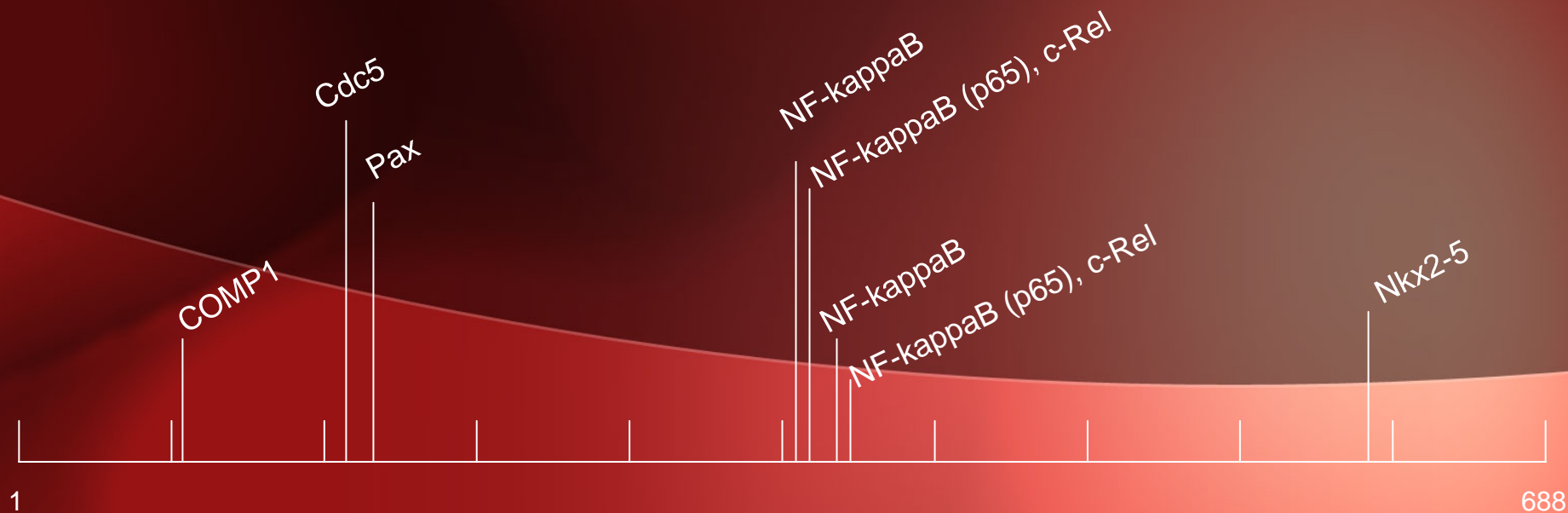
Promoter Regions

- Although we weren't able to analyze HIV-1 with HIV-2, SIV-1, or SIV-2, we compared the Ref Seq promoter region of each species
 - SIV-2 did not have a promoter region defined
- We used MATCH (part of the TRANSFAC database) to predict TFBS within the LTR/U3_R regions of each genome.

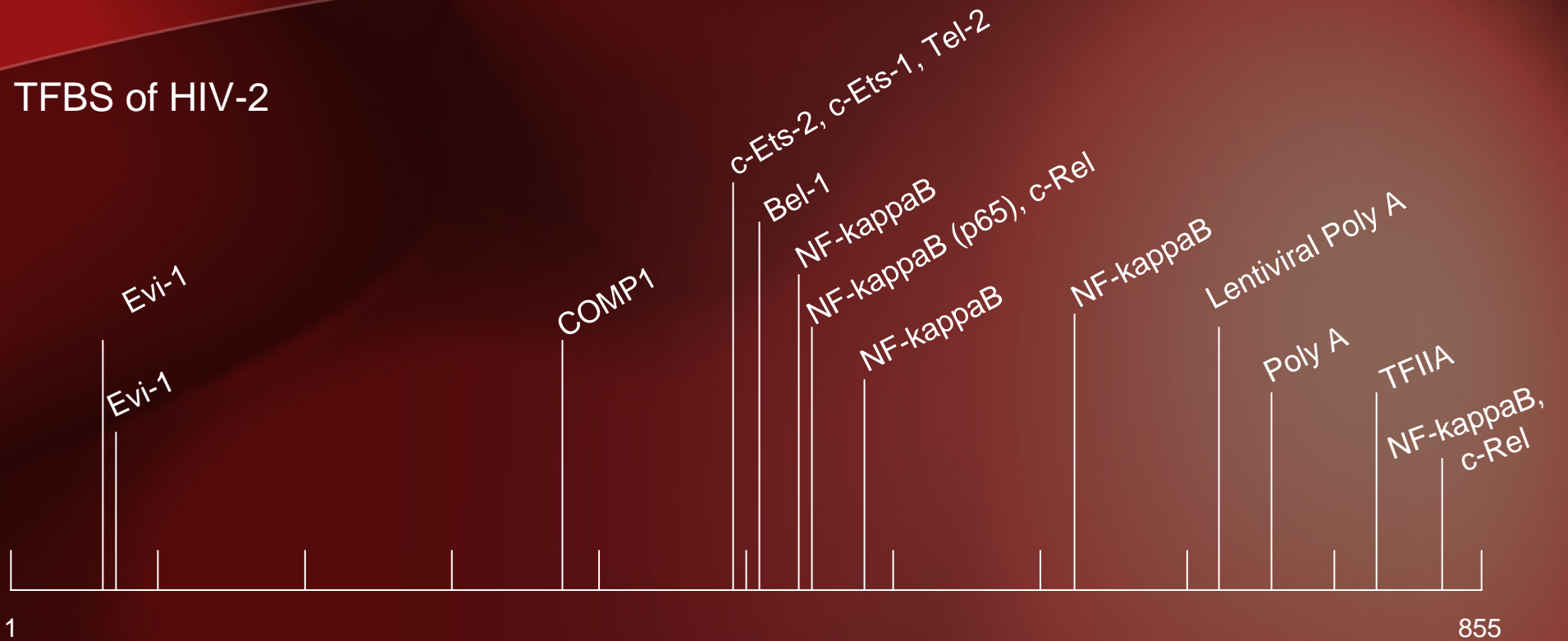
- It uses a scoring matrix generated from known TFBS to predict which TFBS are present within the promoter of the gene
- We predict that there will be many false positives

Results.....

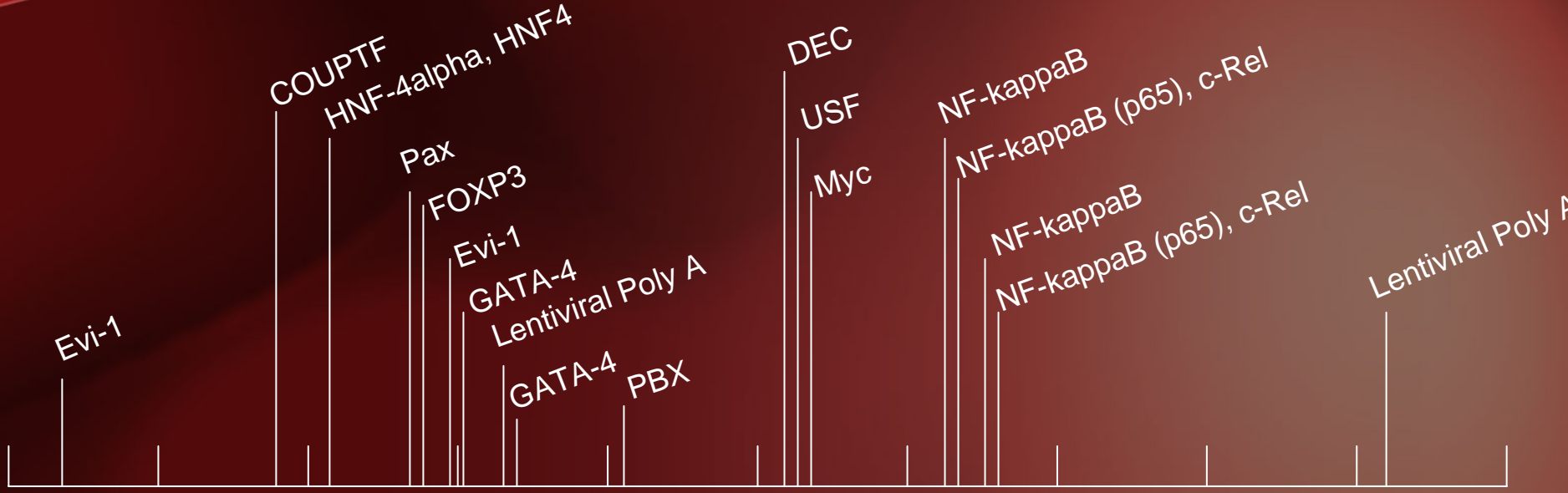
TFBS of SIV-1



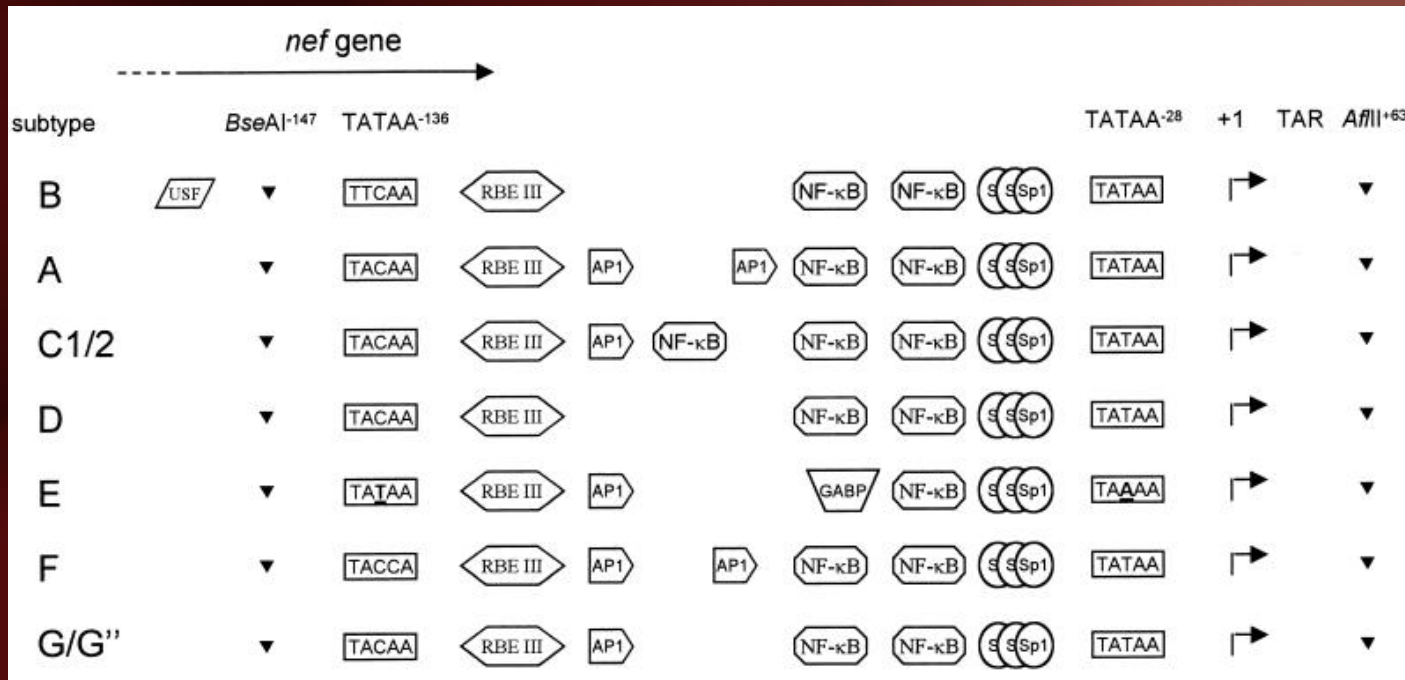
TFBS of HIV-2



TFBS of HIV-1



Comparison to the known TFBS



- Prominence of NF-kappaB site (2)

- We would've liked to see the difference in Θ values across the promoter region. This would confirm and better prove TFBS (the lower Θ values= the more conserved the sequence).

- Use a sliding window of about 100 bp overlapping by 50 bp

- More accurately compare the subtypes of HIV-1
- Further develop the new test used to calculate Θ
- Calculate Θ values for different regions of the promoter region to better prove TFBS

Acknowledgements

- Takis Benos
- David Corcoran
- Shaun Mahony
- Judy Wieber and Rajan Munshi
- Dept. Comp. Bio.- University of Pittsburgh
- All BBSI participants
- NIH-NSF

References

- Kreitman M. and Hudson R. "Inferring the Evolutionary Histories of the *Adh* and *Adh-dup* Loci in *Drosophila melanogaster* From Patterns of Polymorphism and Divergence." Genetics 127 (1991): 565-582.
- McDonald, John H. and Kreitman Martin. "Adaptive protein evolution at the *Adh* locus in *Drosophila*." Nature 351 (1991): 652-654.
- <http://uhavax.hartford.edu/bugl/hiv.htm#types>
- Jeeninga, Rienk E., Hoogenkamp, Maarten, Armand-Ugon, Mercedes, de Baar, Michel, Verhoef, Koen, Berkhout, Ben. "Functional Differences between the Long Terminal Repeat Transcriptional Promoters of Human Immunodeficiency Virus Type 1 Subtypes A through G". Journal of Virology 74:8 (2000): 3740-3751.
- Van Opijnen, Tim, Kamoschinski, Joost, Jeeninga, Rienk E., Berkhout, Ben. "The Human Immunodeficiency Virus Type 1 Promoter Contains a CATA Box instead of a TATA Box for Optimal Transcription and Replication." Journal of Virology 78:13 (2004): 6883-6890.
- UNAIDS (http://www.unaids.org/en/HIV_data/2006GlobalReport/default.asp) presented in: <http://hivinsite.ucsf.edu/global?page=cr-01-00&post=2&cid=US#General%20HIV/AIDS>