**Federal University of Rio de Janeiro**
Institute of Biophysics Carlos Chagas Filho
Bioinformatics for Transcriptomics and Functional Genomics

IBCCF
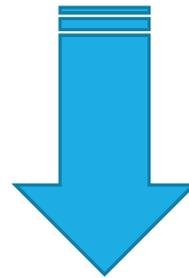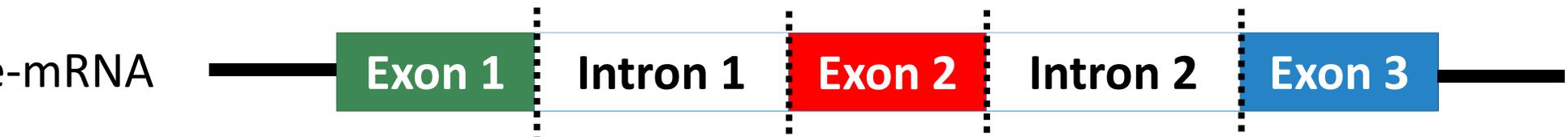
BITFUN

# tomated Functional Annotation of Prot

# Products of Alternatively Spliced Gene

Vitor Coelho & Michael Sammeth

February 2018, Bled Winterseminar
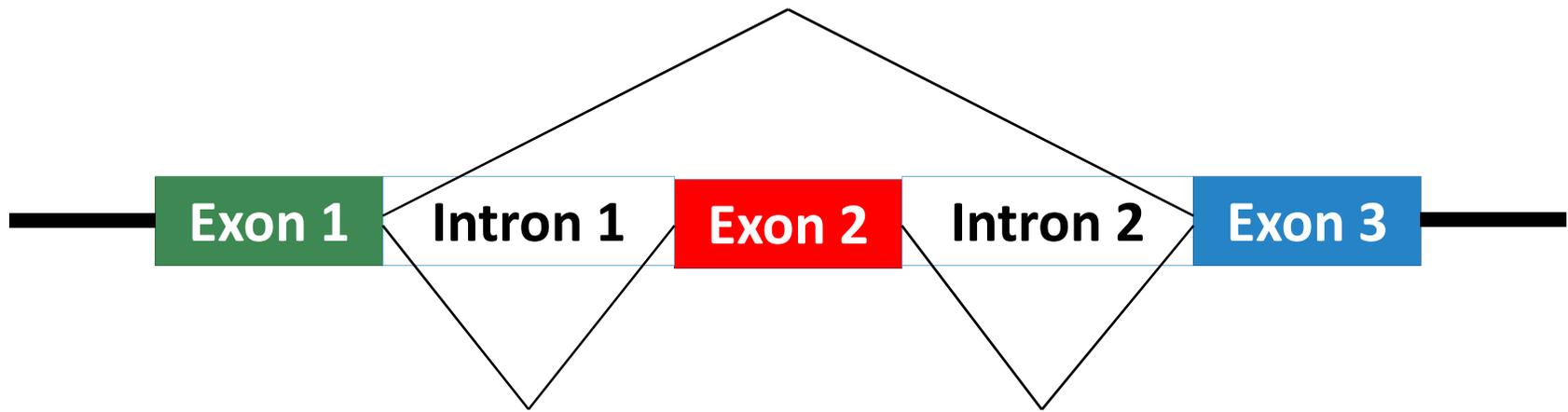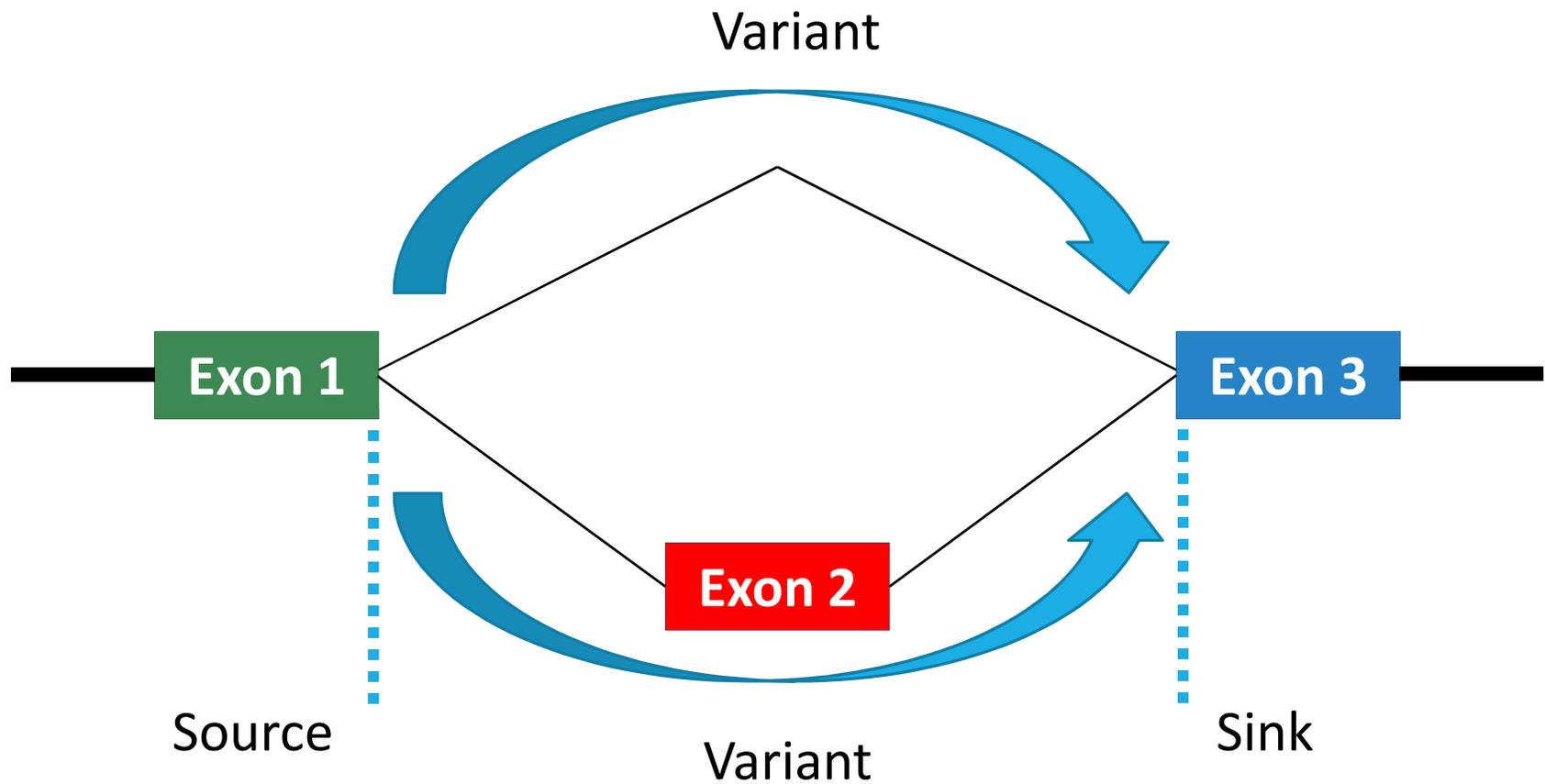
# Alternative Splicing (AS)

# Alternative Splicing Events

# Components of an AS Event



Variant

Exon 1

Exon 3

Exon 2

Source

Variant

Sink

[AStalavista - Foissac&Sammeth 20

# Definition of AS event

o An event e is defined by

$$e = (\, source_e, \, sink_e, \, \mathbb{V}_e)$$

- o Where:
  - *source$_e$*: begin coordinate of event e;
  - *sink$_e$*: end coordinate of event e;
  - $\mathbb{V}_e$ : set of variants;
    - A variant describes an exon/intron structure between *source$_e$* and *sink$_e$* – which can contain multiple transcripts

# Event 1 in TNNT1



$V = \{V_1, V_2, V_3\} \rightarrow$

$V_1 = \{t_1, t_2\}$

$V_2 = \{t_3, t_4\}$

$V_3 = \{t_5\}$

# Event 2 in TNNT1



$$\mathbb{V} = \{V_1, V_2\} \rightarrow$$

$$V_1 = \{t_1, t_2, t_3,$$
$$V_2 = \{t_4\}$$

# Protein Domains and Splicing



Human Protein p85A

Human Protein Lyn

Human Protein Nck1

# The Pfam Database

○ *P*rotein *fam*ilies = database of protein families [Punta et al, 2012]

○ Pfam domains are organized in different sub-groups:

  *Pfam-A*: high quality and manually curated protein families

  *Pfam-B*: automatic predictions

○ protein regions sharing sequence similarity of high significance (*"domains"*) are represented by profile HMMs $\Pi$ $\theta_{\pi}$, each with a specific length $L_\Pi$, a prediction threshold $\theta$, and alignment scores (amino acid scores and gap penalties)

**A** visible layer

$$x_1 \dashrightarrow x_2 \longrightarrow x_{j-1} \longrightarrow x_j \longrightarrow x_{j+1} \dashrightarrow x_{|\mathcal{X}|}$$

**B** hidden layer

$$P_t(M_i|M_{i-1}) \times P_e(x_j|M_i)$$

$M_1 \quad M_{i-1} \quad M_i \quad M_L$

$P_t(D_i|M_{i-1})$

$P_t(M_i|D_{i-1}) \times P_e(x_j|M_i)$

$P_t(M_i|I_{i-1}) \times P_e(x_j|M_i)$

$N \quad B$

$D_1 \quad D_{i-1} \quad D_i \quad D_L$

$E \quad C$

$P_t(D_i|D_{i-1})$

$S$

$T$

$I_1 \quad I_{i-1} \quad I_i$

alignment states:

M – match

D – delete

I – insert

$P_t(I_i|I_i) \times P_e(x_j|I_i)$

$J$

# Pfam Profile HMM file format

```
HMMER3/b [3.0 | March 2010]
NAME  Antimicrobial11
ACC   PF08106.6
DESC  Formaecin family
LENG  16
...
GA    25.00 25.00;
...
```

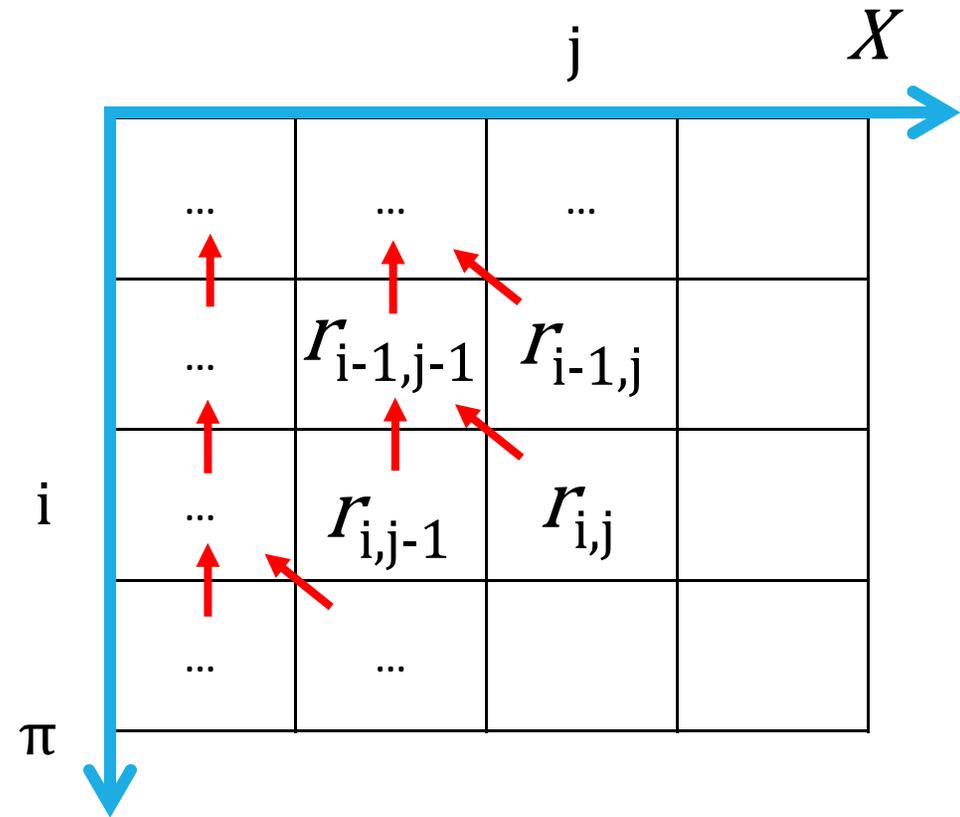| HMM | A | C | D | ... | S | T | V | W | Y | |
|-----|---|---|---|-----|---|---|---|---|---|---|
| | m->m | m->i | m->d | | i->m | i->i | d->m | d->d | | |
| COMPO | 3.68653 | 5.37452 | 4.20831 | ... | 3.83085 | 2.62022 | 2.91489 | 5.77804 | 3.54633 | |
| | 2.68618 | 4.42225 | 2.77519 | ... | 2.37887 | 2.77519 | 2.98518 | 4.58477 | 3.61503 | |
| | 0.01467 | 4.62483 | 5.34718 | | 0.61958 | 0.77255 | 0.00000 | * | | |
| 1 | 3.84125 | 5.47999 | 4.65486 | ... | 4.04240 | 4.37307 | 4.91201 | 6.51694 | 5.82696 | 1 - - |
| | 2.68618 | 4.42225 | 2.77519 | ... | 2.37887 | 2.77519 | 2.98518 | 4.58477 | 3.61503 | |
| | 0.01467 | 4.62483 | 5.34718 | | 0.61958 | 0.77255 | 0.48576 | 0.95510 | | |
| 2 | 4.09153 | 5.75008 | 4.67229 | ... | 4.17197 | 4.33362 | 4.67489 | 6.11461 | 5.14175 | 2 - - |
| | 2.68618 | 4.42225 | 2.77519 | ... | 2.37887 | 2.77519 | 2.98518 | 4.58477 | 3.61503 | |
| | 0.01467 | 4.62483 | 5.34718 | | 0.61958 | 0.77255 | 0.48576 | 0.95510 | | |
| ... | | | | | | | | | | |
| 15 | 4.09153 | 5.75008 | 4.67229 | ... | 4.17197 | 4.33362 | 4.67489 | 6.11461 | 5.14175 | 15 - - |
| | 2.68618 | 4.42225 | 2.77519 | ... | 2.37887 | 2.77519 | 2.98518 | 4.58477 | 3.61503 | |
| | 0.01467 | 4.62483 | 5.34718 | | 0.61958 | 0.77255 | 0.48576 | 0.95510 | | |
| 16 | 4.12723 | 5.39816 | 5.40081 | ... | 4.69892 | 4.40337 | 3.12238 | 5.70016 | 4.58094 | 16 - - |
| | 2.68618 | 4.42225 | 2.77519 | ... | 2.37887 | 2.77519 | 2.98518 | 4.58477 | 3.61503 | |
| | 0.00990 | 4.62006 | * | | 0.61958 | 0.77255 | 0.00000 | * | | |

# Viterbi algorithm

o Dynamic Programming
- recursive algorithm
- divides exponential problem into polynomial sub-problems

o Calculate best alignment score between prefixes:
- Profile HMM $\pi[1; i]$ and query sequence $X[1; j]$



o path of best alignment is stored by pointers to the cell that led to highest score

# Viterbi algorithm

$$r_{i,j}^{M} = \log_2\left(\frac{P_e(x_j|M_i)}{P_0(x_j)}\right) + \max \begin{cases} r_{i-1,j-1}^{M} & +\log_2(P_t(M_i|M_{i-1})) \\ r_{i-1,j-1}^{I} & +\log_2(P_t(M_i|I_{i-1})) \\ r_{i-1,j-1}^{D} & +\log_2(P_t(M_i|D_{i-1})) \end{cases}$$

$$r_{i,j}^{I} = \log_2\left(\frac{P_e(x_j|I_i)}{P_0(x_j)}\right) + \max \begin{cases} r_{i,j-1}^{M} & +\log_2(P_t(I_i|M_i)) \\ r_{i,j-1}^{I} & +\log_2(P_t(I_i|I_i)) \end{cases}$$

$$r_{i,j}^{D} = \max \begin{cases} r_{i,j-1}^{M} & +\log_2(P_t(D_i|M_{i-1})) \\ r_{i,j-1}^{D} & +\log_2(P_t(D_i|D_{i-1})) \end{cases}$$

**vcoelho2**   está errado as probabilidades de emissao... estao como probabilidades de transicao
vitorlc; 22/2/2015

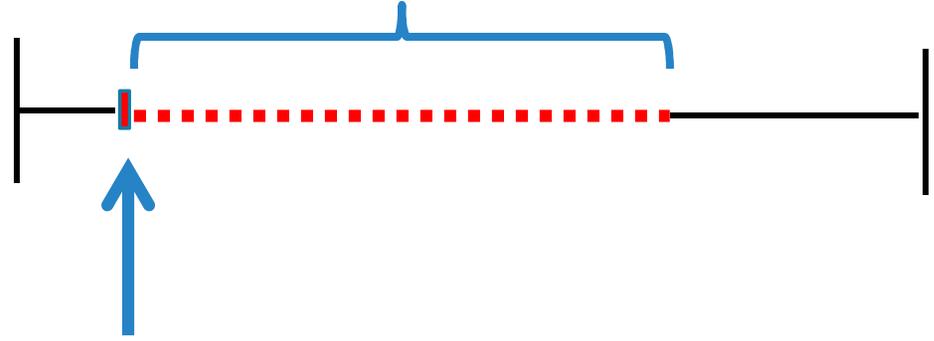# Aim: Predict Domain Alterations by AS

# Where to scan for AS domains?

$$\Delta_\pi = ?$$

*gaps*     *3L-1*

$$\Delta_\pi = ?$$

*3L-1*   *gaps*

$e_1$

# How long has to be $\Delta_\Pi$?

# What is the minimum required $\Delta_\Pi$?

$$\Delta_\pi = \left\lceil \frac{(L_\pi \times 3) - 1}{3} + \frac{\omega_\pi(0) - \theta_\pi - \alpha_\pi}{\beta_\pi} \right\rceil \times 3$$

$L$:   *length* of profile HMM π

$\omega_\pi(1)$:   *optimal alignment* score for π (optimal *suffix alignment* starting at position

$\theta_\pi$:   Domain *gathering threshold* (Pfam determined) for relevant domain score

$\alpha_\pi$:   absolute value of max. bit score to ***open*** *an insertion* ("gap", state I of the mo

$\beta_\pi$:   absolute value of max. bit score to ***extend*** *an insertion*

# Extension of AS Event

$\Delta_\pi$

vent

# Fusing Events

# Splitting Variants



$$\left\{\begin{array}{l} V_1 = \{t_1, t_2\} \\ V_2 = \{t_3, t_4\} \\ V_3 = \{t_5\} \end{array}\right.$$

$e_2$

$$\mathbb{V}e_2e_3 = \left\{\begin{array}{l} V_1 = \{t_1, t_2\} \\ V_2 = \{t_3\} \\ V_3 = \{t_4\} \\ V_4 = \{t_5\} \end{array}\right.$$

$e_3$

$$\mathbb{V}e_3 = \left\{\begin{array}{l} V_1 = \{t_1, t_2, t_3, t_5\} \\ V_2 = \{t_4\} \end{array}\right.$$

# Algorithm: Split Variants

**Data**: $\{\mathbb{V}_{e_1}, \ldots, \mathbb{V}_{e_n}\}$ = variant sets

**Result**: $\mathbb{V}_{e_1,\ldots,e_n}$ = variant set of fused event

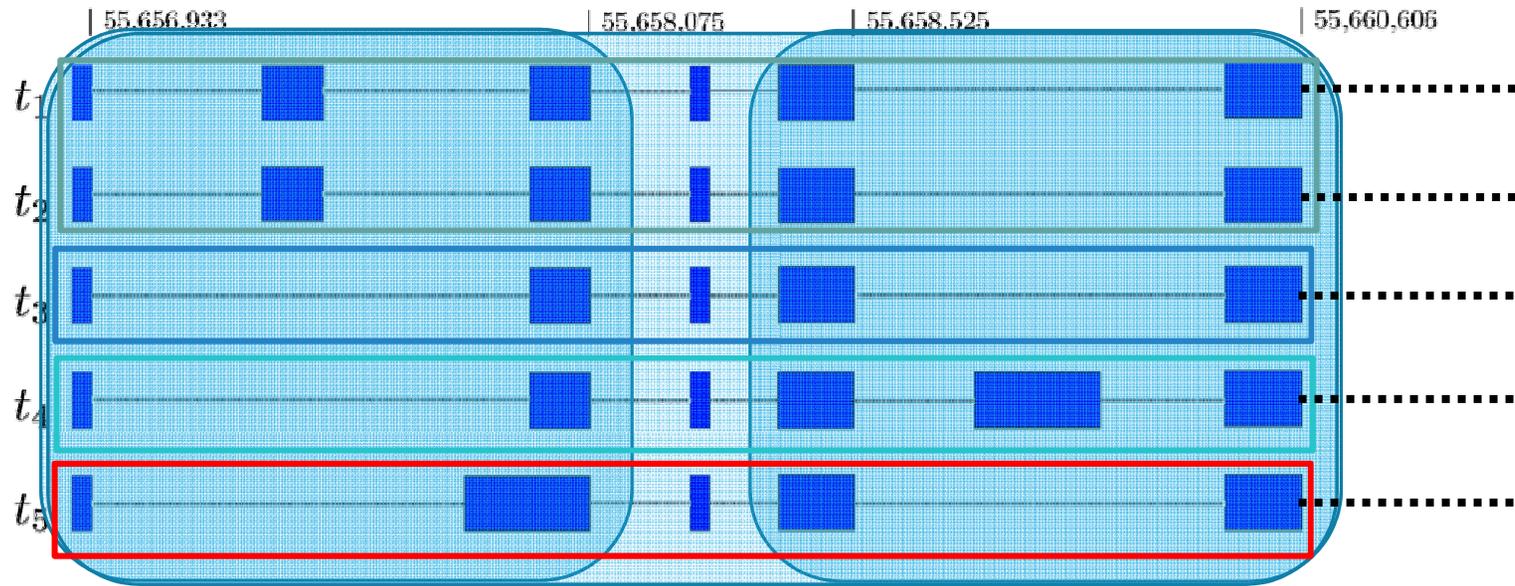1   $V^1_{e_1,\ldots,e_n} \leftarrow \bigcup\limits_{i=1}^{n} (V_{e_i} \in \mathbb{V}_{e_i})$

2   $m \leftarrow |V^1_{e_1,\ldots,e_n}|$

3   $\mathbb{V}_{e_1,\ldots,e_n} \leftarrow V^1_{e_1,\ldots,e_n}$

4   **for** $x \leftarrow 1$ **to** $m$ **do**

5     $W \leftarrow \emptyset$

6     **for** $y \leftarrow x + 1$ **to** $m$ **do**

7       **for** $i \leftarrow 1$ **to** $n$ **do**

8         **if** $(t_x \in V^j_{e_i}) \wedge (t_y \in V^k_{e_i}) \wedge (t_x \in V^l_{e_1,\ldots,e_n}) \wedge (t_y \in V^l_{e_1,\ldots,e_n})$ **then**

9           $V^l_{e_1,\ldots,e_n} \leftarrow V^l_{e_1,\ldots,e_n} \setminus \{t_y\}$

10          $W \leftarrow W \cup \{t_y\}$

11     $V^{|\mathbb{V}_{e_1,\ldots,e_n}|+1}_{e_1,\ldots,e_n} \leftarrow W;$

      $\mathbb{V}_{e_1,\ldots,e_n} \leftarrow \mathbb{V}_{e_1,\ldots,e_n} \cup V^{|\mathbb{V}_{e_1,\ldots,e_n}|+1}_{e_1,\ldots,e_n}$

# Optimization: Branch-and-bound Condition for the Viterbi Algorithm

o improve alignment performance for finding entire domains (no partial hits).
  - reduce cells of DP matrix to calculate pruning sub-solutions that can no longer produce a relevant alignment (score $\geq \theta_{\pi}$).

o condition for the feasibility of a sub-solution

$$ r_{i,j}^{z} + \omega(i) - \log_2(|X|) \geq \theta_{\pi} $$

omega = optimal suffix alignment  [i..|X|]

$\log_2(|X|)$ = score normalization factor

| | Flybase (dm3) | Wormbase (ce6) | RefSeq (hg19) | UCSC (hg19) | GEN (h |
|---|---|---|---|---|---|
| es | 3,541 | 4,124 | 10,710 | 15,091 | |
| es without (valid) domains | 503 | 795 | 1,523 | 1,899 | |
| ed AS genes | 3,038 | 3,329 | 9,187 | 13,192 | |
| ed AS transcripts | 9,779 | 11,836 | 33,290 | 61,525 | 1 |
| redicted domains | 777 | 1,497 | 1,956 | 3,078 | |
| events / AS gene | 2.3 | 2.7 | 2.9 | 4.7 | |
| tive search [min] | 12.9 | 35.9 | 688.5 | 716.0 | |
| UNK search [min] | 4.02 | 7.9 | 29.9 | 37.6 | |
| **peedup Gain** | **68.84%** | **78%** | **95.7%** | **94.7%** | |

# AS Density as a determinant of Fused Events (overhead)

# Runtime as a Function of the Density/Fusion of AS Events

$\Delta_\pi$          $\Delta_\pi$          $\Delta_\pi$          $\Delta_\pi$

$e_1$                $e_2$

# Runtime as a Function of the Density of AS Events

# Examples

## Pattern 0_1d2a

**Dimension:**

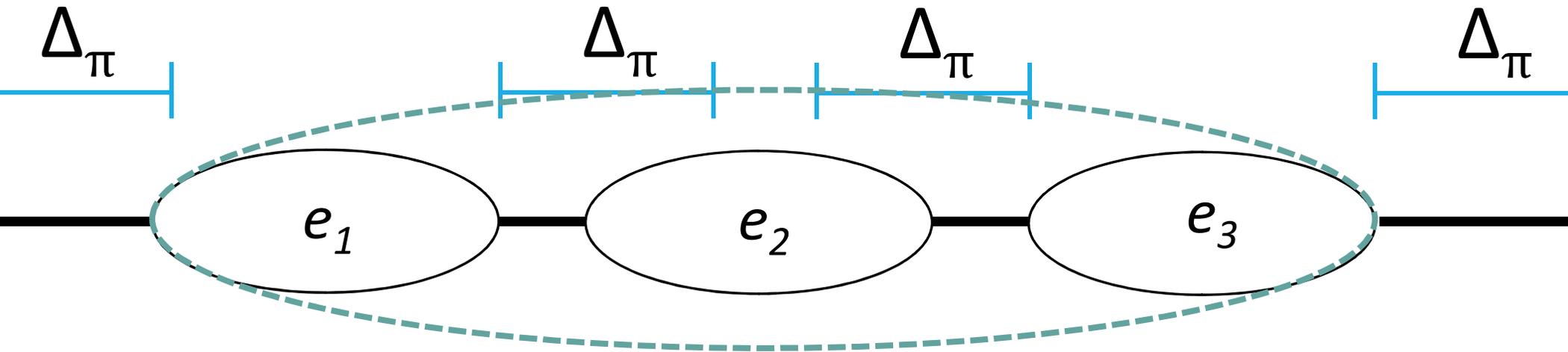| Locus | Domain | Scores | Schema | AS code | Variants |
|---|---|---|---|---|---|
| chr3L:13477541 | PF02179.15 | 72.51<br>72.16 | | 0<br>1^2- | FBtr0075808/FBtr0113425<br>FBtr0075806/FBtr0075807 |
| chr3R:2260306 | PF01369.19 | 116.92<br>117.47 | | 0<br>1^2- | FBtr0113396/FBtr0334559/FBt<br>FBtr0113395/FBtr0302187 |
| chrX:2344559-2 | PF00238.18 | 43.84<br>43.84 | | 0<br>1^2- | FBtr0070405<br>FBtr0333730 |
| chr3L:9728149- | PF08445.9 | 22.50<br>N/A | | 0<br>1^2- | FBtr0089411<br>FBtr0089410/FBtr0114566 |
| chr3R:3002284 | PF00615.18 | N/A<br>68.82 | | 0<br>1^2- | FBtr0300600<br>FBtr0085553/FBtr0085554/FBt |
| chrX:7933774-7 | PF00411.18 | 164.54<br>164.54 | | 0<br>1^2- | FBtr0071095<br>FBtr0071094 |
| chrX:1471348-1 | PF13903.5 | N/A<br>76.91 | | 0<br>1^2- | FBtr0100198<br>FBtr0070255/FBtr0070256 |
| chr2R:2242625 | PF01529.19 | 76.68<br>64.88 | | 0<br>1^2- | FBtr0071841<br>FBtr0071840 |
| chr3L:1608067- | PF16870.4 | 148.32<br>183.83 | | 0<br>1^2- | FBtr0091793/FBtr0091795/FBt<br>FBtr0091792/FBtr0091794/FBt |

# Pattern 0_1a2d3a4d

**Dimension:**

| Locus | Domain | Scores | Schema | AS code | Variants |
|-------|--------|--------|--------|---------|----------|
| chr3R:152813 7: | PF00261.19 | N/A 331.06 | | 0 1-2^3-4^ | FBtr0089957/FBtr0089959/FBt FBtr0089965 |
| chr3R:152813 7: | N/A | N/A N/A | | 0 1-2^3-4^ | FBtr0089957/FBtr0089959/FBt FBtr0089965 |
| chr3L:213187 74 | PF00071.21 | N/A 193.43 | | 0 1-2^3-4^ | FBtr0112677/FBtr0345428 FBtr0112678/FBtr0331794 |
| chr2R:185229 7: | PF00209.17 | N/A 595.43 | | 0 1-2^3-4^ | FBtr0304867 FBtr0304865 |
| chr3L:213187 74 | PF00025.20 | N/A 36.67 | | 0 1-2^3-4^ | FBtr0112677/FBtr0345428 FBtr0112678/FBtr0331794 |
| chr3R:152813 7: | N/A | N/A N/A | | 0 1-2^3-4^ | FBtr0089957/FBtr0089959/FBt FBtr0089965 |
| chr3R:152813 7: | PF12718.6 | N/A 54.11 | | 0 1-2^3-4^ | FBtr0089957/FBtr0089959/FBt FBtr0089965 |
| chr2L:149839 5( | PF10204.8 | N/A 204.73 | | 0 1-2^3-4^ | FBtr0332521 FBtr0080676/FBtr0111014/FBt |
| chrX:17175347- | PF12053.7 | N/A 28.60 | | 0 1-2^3-4^ | FBtr0074388/FBtr0343763 FBtr0111000/FBtr0343764 |
| chr3L:213187 74 | PF08477.12 | N/A 107.58 | | 0 1-2^3-4^ | FBtr0112677/FBtr0345428 FBtr0112678/FBtr0331794 |

# Pattern 1a2e_3a4e

Dimension:

| Locus | Domain | Scores | Schema | AS code | Variants |
|-------|--------|--------|--------|---------|----------|
| chr2L:5762553- | PF01607.23 | 39.92<br>31.03 | | 1-2]<br>3-4] | FBtr0079163<br>FBtr0079164 |
| chr3R:6088799- | PF01607.23 | N/A<br>31.75 | | 1-2]<br>3-4] | FBtr0078619<br>FBtr0112922 |
| chr3R:2241454 | PF10193.8 | 91.15<br>N/A | | 1-2]<br>3-4] | FBtr0299516<br>FBtr0299515 |
| chr3L:24537911 | PF00333.19 | N/A<br>58.84 | | 1-2]<br>3-4] | FBtr0111144<br>FBtr0111146/FBtr0111147 |
| chr3R:19754254 | PF02932.15 | 84.76<br>77.53 | | 1-2]<br>3-4] | FBtr0335417/FBtr0335418/FBtr03<br>FBtr0335416/FBtr0335419/FBtr03 |
| chr3R:22458321 | PF04500.15 | 44.91<br>76.18 | | 1-2]<br>3-4] | FBtr0112608<br>FBtr0112609 |
| chr3R:13222951 | PF00135.27 | N/A<br>520.36 | | 1-2]<br>3-4] | FBtr0082780<br>FBtr0335223 |
| chr3R:9726020- | PF00250.17 | 96.10<br>N/A | | 1-2]<br>3-4] | FBtr0329937<br>FBtr0304869 |
| chr2R:9400531- | PF01569.20 | 105.50<br>78.20 | | 1-2]<br>3-4] | FBtr0306111<br>FBtr0088552 |
| chr3L:24537911 | PF03719.14 | N/A<br>75.70 | | 1-2]<br>3-4] | FBtr0111144<br>FBtr0111146/FBtr0111147 |

# Pattern degree-6_dimension-2

**Dimension:**

| Locus | Domain | Scores | Schema | AS code | Variants |
|---|---|---|---|---|---|
| chr3L:16860622 | PF01529.19 | 95.41<br>N/A | | 1[2^3-4^6-<br>5[ | FBtr0075336<br>FBtr0332702 |
| chr2L:3713360- | PF00686.18 | 31.28<br>N/A | | 1[2^3-5^<br>4[6^ | FBtr0077537<br>FBtr0077538 |
| chr2L:297880-3 | PF00017.23 | 69.97<br>N/A | | 1[2^3-4^<br>5[6^ | FBtr0331205/FBtr0331208<br>FBtr0331207 |
| chr3R:1916009 | PF01189.16 | 103.56<br>N/A | | 1^3-4^5-6]<br>2] | FBtr0299957<br>FBtr0299958 |
| chrX:3877628-3 | PF00854.20 | 308.74<br>312.10 | | 1-2]<br>3-4^5-6] | FBtr0070609/FBtr0070610/FBtr<br>FBtr0070608 |
| chr3L:1887004 | PF04139.12 | 56.20<br>56.20 | | 1^3-4^5-6]<br>2] | FBtr0075080<br>FBtr0112917 |
| chr2L:297880-3 | PF16454.4 | 164.68<br>N/A | | 1[2^3-4^<br>5[6^ | FBtr0331205/FBtr0331208<br>FBtr0331207 |
| chr3L:6704860- | PF00071.21 | 100.45<br>N/A | | 1[2^<br>3[4^5-6^ | FBtr0076942<br>FBtr0333072 |
| chr3L:1162354 | PF00248.20 | 175.20<br>178.48 | | 1[2^3-4^<br>5[6^ | FBtr0076138/FBtr0331546<br>FBtr0076139 |
| chr3R:1825352 | PF04433.16 | N/A<br>32.17 | | 0<br>1-2^3-4^5-6^ | FBtr0310379<br>FBtr0310384 |
| chr3R:9365398- | PF00291.24 | 268.32<br>N/A | | 1[2^3-4^<br>5[6^ | FBtr0082026<br>FBtr0082027 |

# Pattern degree-6_dimension-3

**Dimension:**

| Locus | Domain | Scores | Schema | AS code | Variants |
|---|---|---|---|---|---|
| chr3L:2127619⁵ | PF00501.27 | 297.96 297.96 N/A |  | 1[3^ 2[3^ 4[5^ | FBtr0078413 FBtr0333126 FBtr0078414 |
| chrX:769945-7⁷ | PF01392.21 | 92.34 92.34 45.05 |  | 1[3^ 2[3^ 4[5^ | FBtr0345511 FBtr0070138 FBtr0112928 |
| chrX:4081354-⁴ | PF01530.17 | N/A 52.38 52.50 |  | 0 1-2^3-4^ 3-4^ | FBtr0310491/FBtr0310494/FBtr FBtr0310493 FBtr0310492/FBtr0310495 |
| chr2L:12434452 | PF07145.14 | 26.32 26.32 N/A |  | 1[3^ 2[3^ 4[5^ | FBtr0080369/FBtr0333089 FBtr0344859 FBtr0333090 |
| chr2L:22240154 | PF01656.22 | 34.77 34.77 N/A |  | 1-4] 1-5] 2-3] | FBtr0299875 FBtr0334168 FBtr0334169 |
| chr2R:1010974( | PF00365.19 | 368.27 377.77 375.62 |  | 1-2^ 3-4^ 5-6^ | FBtr0304925/FBtr0333067 FBtr0088420/FBtr0088422 FBtr0088421 |
| chr3R:7861139- | PF15501.5 | N/A 244.27 132.95 |  | 0 1-2^ 1-2^3-4^ | FBtr0333828 FBtr0333827 FBtr0081770/FBtr0113202/FBtr |
| chr3L:17539727 | PF00168.29 | 80.92 79.84 80.15 |  | 1^4- 2^3- 2^4- | FBtr0299642/FBtr0331844 FBtr0331843 FBtr0299645/FBtr0299646/FBtr |

# AstaFunk

**http://astafunk.sammeth.net**

**Vitor Lima Coelho & Michael Sammeth**

**micha@sammeth.net**

**vitorcoelho@biof.ufrj.br**

CNPq
elho Nacional de Desenvolvimento
fico e Tecnológico