

Trip into the Desert

Drosophila melanogaster & Co

Sven Findeiß

Lehrstuhl für Bioinformatik
University of Leipzig

Outline

Introduction

What I've seen

Motivation

Focus of my work

Blat search for orthologous protein sequences

Outline

Introduction

What I've seen

Motivation

Focus of my work

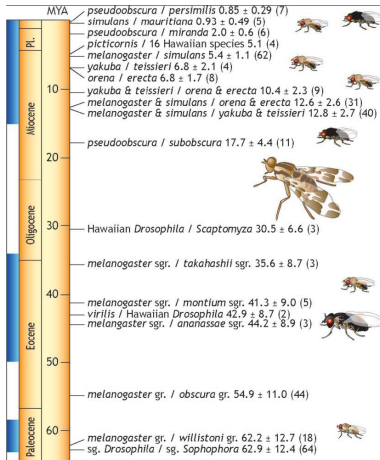
Blat search for orthologous protein sequences





Center for Evolutionary Functional Genomics (EFG)

Current research achievements:



Outline

Introduction

What I've seen

Motivation

Focus of my work

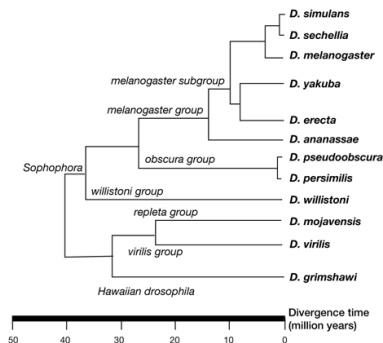
Blat search for orthologous protein sequences

Motivation

Why analyze *Drosophila* genomes?

- *Drosophila melanogaster* is a well studied model organism
- high homology of *Drosophila melanogaster* and human (584 *Drosophila* genes related to human disease ^{a)})
- 11 recently sequenced non-*melanogaster* species

^{a)}<http://superfly.ucsd.edu/homophila/>



Outline

Introduction

What I've seen

Motivation

Focus of my work

Blat search for orthologous protein sequences

Idea

Use the blat algorithm to search for orthologous protein sequences in the 11 non-*melanogaster* species.

- *Drosophila melanogaster* is a well studied model organism
⇒ list of annotated protein sequences downloaded from
`http://genome.ucsc.edu`
- *melanogaster* and recently sequenced genomes
downloaded from
`http://rana.lbl.gov/drosophila/caf1.html`
- *Drosophila melanogaster* was used as reference sequence

Some details

Calculation of Scores

$$UCSC_{score} = 3 \times (match + (repMatch \gg 1)) \\ - 3 \times misMatch - qNumInsert - tNumInsert$$

$$INTRA_{score} = \frac{best UCSC_{score} - second\ best\ UCSC_{score}}{best\ UCSC_{score}}$$

⇒ value close to 1 shows the uniqueness of a hit

$$INTER_{score} = \frac{best\ UCSC_{score}\ current}{best\ UCSC_{score}\ reference}$$

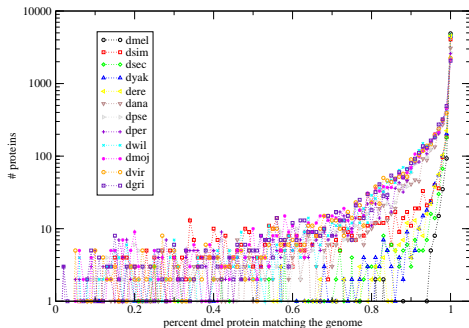
⇒ value close to 1 means high similarity to reference hit
(dmel protein against dmel genome)



The output (1)

emacs@localhost.localdomain										
File Edit Options Buffers Tools Help										
#NAME	CHROMOSOM	PROTEIN	STRAND	SCORE	INTRA	INTER	NTERM	CTERM	HITSTART	HITEND
dnel NP_001007096.1	dnel_3L	NP_001007096.1	-	4557	0.941409	1.000000	0.000000	0.000000	19685705	19692073
dana NP_001007096.1	dana_scaffold_13337	NP_001007096.1	+	3755	0.934221	0.824007	0.000000	0.000000	17177309	17183397
dere NP_001007096.1	dere_scaffold_4784	NP_001007096.1	-	4399	0.880200	0.965328	0.000000	0.000000	19506009	19512662
dgrj NP_001007096.1	dgrj_scaffold_15110	NP_001007096.1	+	3114	0.885035	0.683344	0.000000	0.000000	16662632	16669501
dnoj NP_001007096.1	dnoj_scaffold_6680	NP_001007096.1	-	3307	0.888419	0.725697	0.000000	0.000000	12260037	12266315
dper NP_001007096.1	dper_super_48	NP_001007096.1	-	2743	0.875319	0.601931	0.000000	0.276721	17526	574971
dpse NP_001007096.1	dpse_ChIR_group8	NP_001007096.1	-	3747	0.916200	0.822251	0.000000	0.004590	7788295	7794126
dsec NP_001007096.1	dsec_super_29	NP_001007096.1	-	4366	0.948923	0.958086	0.000000	0.000000	411269	417574
dsin NP_001007096.1	dsin_chr3L	NP_001007096.1	-	4511	0.920195	0.999906	0.000000	0.000000	19039455	19046110
dvir NP_001007096.1	dvir_scaffold_13049	NP_001007096.1	+	3380	0.900000	0.741716	0.000000	0.000000	18965491	18972621
dwl NP_001007096.1	dwl_scaffold_180949	NP_001007096.1	-	3345	0.889088	0.734036	0.000000	0.000000	4854927	4861693
dyak NP_001007096.1	dyak_chr3L	NP_001007096.1	+	4403	0.399727	0.966206	0.000000	0.000000	20807162	20813839
dnel NP_001007097.1	dnel_3L	NP_001007097.1	-	4557	0.941409	1.000000	0.000000	0.000000	19685705	19692073
dana NP_001007097.1	dana_scaffold_13337	NP_001007097.1	+	3755	0.934221	0.824007	0.000000	0.000000	17177309	17183397
dere NP_001007097.1	dere_scaffold_4784	NP_001007097.1	-	4399	0.880200	0.965328	0.000000	0.000000	19506009	19512662
dgrj NP_001007097.1	dgrj_scaffold_15110	NP_001007097.1	+	3114	0.885035	0.683344	0.000000	0.000000	16662632	16669501
dnoj NP_001007097.1	dnoj_scaffold_6680	NP_001007097.1	-	3307	0.88419	0.725697	0.000000	0.000000	12260037	12266315
dper NP_001007097.1	dper_super_48	NP_001007097.1	-	2743	0.875319	0.601931	0.000000	0.276721	17526	574971
dpse NP_001007097.1	dpse_ChIR_group8	NP_001007097.1	-	3747	0.916200	0.822251	0.000000	0.004590	7788295	7794126
dsec NP_001007097.1	dsec_super_29	NP_001007097.1	-	4366	0.948923	0.958086	0.000000	0.000000	411269	417574
dsin NP_001007097.1	dsin_chr3L	NP_001007097.1	-	4511	0.920195	0.999906	0.000000	0.000000	19039455	19046110
dvir NP_001007097.1	dvir_scaffold_13049	NP_001007097.1	+	3380	0.900000	0.741716	0.000000	0.000000	18965491	18972621
dwl NP_001007097.1	dwl_scaffold_180949	NP_001007097.1	-	3345	0.889088	0.734036	0.000000	0.000000	4854927	4861693
dyak NP_001007097.1	dyak_chr3L	NP_001007097.1	+	4403	0.399727	0.965206	0.000000	0.000000	20907162	20913839
dnel NP_001014582.1	dnel_3L	NP_001014582.1	-	2006	0.865902	1.000000	0.000000	0.000000	12437218	12441565
dana NP_001014582.1	dana_scaffold_13337	NP_001014582.1	+	1499	0.723149	0.747258	0.000000	0.001493	1201915	1205972
dere NP_001014582.1	dere_scaffold_4784	NP_001014582.1	-	1901	0.885324	0.947657	0.000000	0.000000	12452669	12457014
dgrj NP_001014582.1	dgrj_scaffold_15110	NP_001014582.1	+	1224	0.804739	0.610169	0.000000	0.019403	5059343	6158280
dnoj NP_001014582.1	dnoj_scaffold_6680	NP_001014582.1	-	1219	0.812961	0.607677	0.000000	0.000000	6057840	6061120
dper NP_001014582.1	dper_super_9	NP_001014582.1	-	1529	0.661871	0.762213	0.000000	0.000000	1537584	1541222
dpse NP_001014582.1	dpse_ChIR_group6	NP_001014582.1	+	1435	0.397362	0.965206	0.000000	0.000000	3245717	3245819
dsec NP_001014582.1	dsec_super_0	NP_001014582.1	-	1926	0.868640	0.960120	0.000000	0.000000	4640331	4644534
dsin NP_001014582.1	dsin_chr3L	NP_001014582.1	-	1929	0.881286	0.961615	0.000000	0.000000	11838230	11842551
dvir NP_001014582.1	dvir_scaffold_13049	NP_001014582.1	+	1213	0.810387	0.604686	0.000000	0.000000	9289088	9292745

----- dnel.output (Fundamental)-----L1--C0--Top-----
 Long Line Truncation enabled



Conclusion:

- fraction of dmel genes matching the genomes represent the relationship of flies (exception dsim)
- data quality: most of the dmel proteins match other drosophila genomes with more than 80% of the sequence