

Some thoughts on designing a genomic query language

Limsoon Wong

Short talk for GeCo Workshop, Como, March 2019



NUS
National University
of Singapore

National University of Singapore

Two perspectives on a query language

Surface syntax

Easy to read, understand, & write queries correctly

Sufficient power to express needed **queries**

Prevent expensive queries

```
{ {x, z) | (x,y) ∈ A, (y',z) ∈ B, y = y' }
```

```
select (a.x, b.z)
from a in A, b in B
where a.y = b.y'
```

Abstract syntax

Easy to analyze, manipulate, and optimize

Easy to cater for extensions

Sufficient power to express needed **algorithms**

```
U{ U{ if a.y = b.y'
      then { (a.x, bz) } else {}
    | b ∈ B }
  | a ∈ A }
```

Compositionality & orthogonality are key principles for query language design

Compositionality & orthogonality in NRC

Types

$$s, t ::= \mid \text{bool} \mid b \mid s \times t \mid \{s\}$$

Expressions, *constructs* are provided for each type orthogonally

$$\frac{}{x^s : s} \quad \frac{e_1 : s \quad e_2 : t}{(e_1, e_2) : s \times t} \quad \frac{e : s \times t}{\pi_1 e : s \quad \pi_2 e : t}$$

$$\frac{}{\text{true} : \text{bool}} \quad \frac{}{\text{false} : \text{bool}} \quad \frac{e_1 : \text{bool} \quad e_2 : s \quad e_3 : s}{\text{if } e_1 \text{ then } e_2 \text{ else } e_3 : s}$$

$$\frac{}{\{ \}^s : \{s\}} \quad \frac{e : s}{\{e\} : \{s\}} \quad \frac{e_1 : \{s\} \quad e_2 : \{s\}}{e_1 \cup e_2 : \{s\}}$$

$$\frac{e_1 : \{s\} \quad e_2 : \{t\}}{\cup \{e_1 \mid x^t \in e_2\} : \{s\}} \quad \frac{e : \{s\}}{\text{empty } e : \text{bool}} \quad \frac{e_1 : s \quad e_2 : s}{e_1 = e_2 : \text{bool}}$$

$\cup \{e_1 \mid x \in e_2\}$ means
 $f(o_1) \cup \dots \cup f(o_n)$,
 where $f(x) = e_1$
 and $\{o_1, \dots, o_n\} = e_2$

Translating into comprehension syntax

$$\cup \{! e_1 \mid x \in e_2\} = \{! y \mid x \in e_2, y \in e_1\}$$

Translating from comprehension syntax

$$\{! e_1 \mid x \in e_2, \Delta\} = \cup \{! e_1 \mid \Delta\} \mid x \in e_2\}$$

$$\{! e_1 \mid C, \Delta\} = \text{if } C \text{ then } \{! e_1 \mid \Delta\} \text{ else } \{\}$$

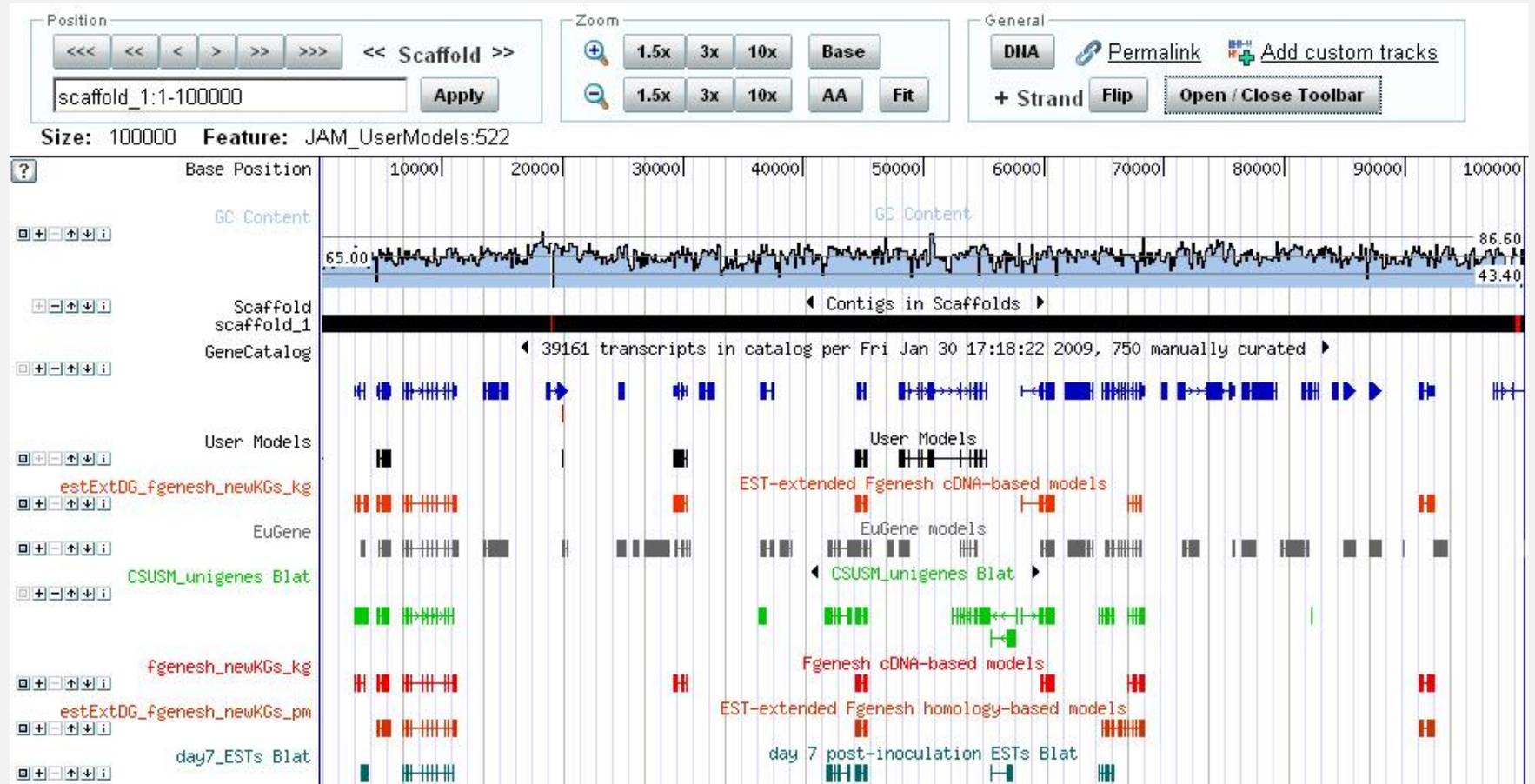
$$\{! e_1 \mid \} = \{! e_1\}$$

⇒ Treat comprehension as a nice syntactic sugar

Genomic data

Loci Tracks Annotations

E.g. you see a row denoted “Base position”; this is the coordinates on a reference genome. The rest of the rows (e.g. “Gene catalog” and “day7_ESTs blat”) are “tracks” bearing different kinds of annotations. Each track corresponds to one kind of experiments, one kind of predictions, etc.; e.g. “day7_ESTs blat” are short stretches of RNAs from a day-7 sample that have been mapped (using a tool called “blat”) to specific positions on the reference genome.



Genomic data types

An **annotation** datatype `!t` and its subtype **landmark** `!!t`

of type `t = (#name: string, #pval: real, ...)`

are represented as `(#loc: Loc, #anno: (name: string, #pval: real, ...))`

A **track** datatype `{!t}` and its subtype **landmark track** `{!!t}`

are represented as `{ !t }` and `{ !!t }`



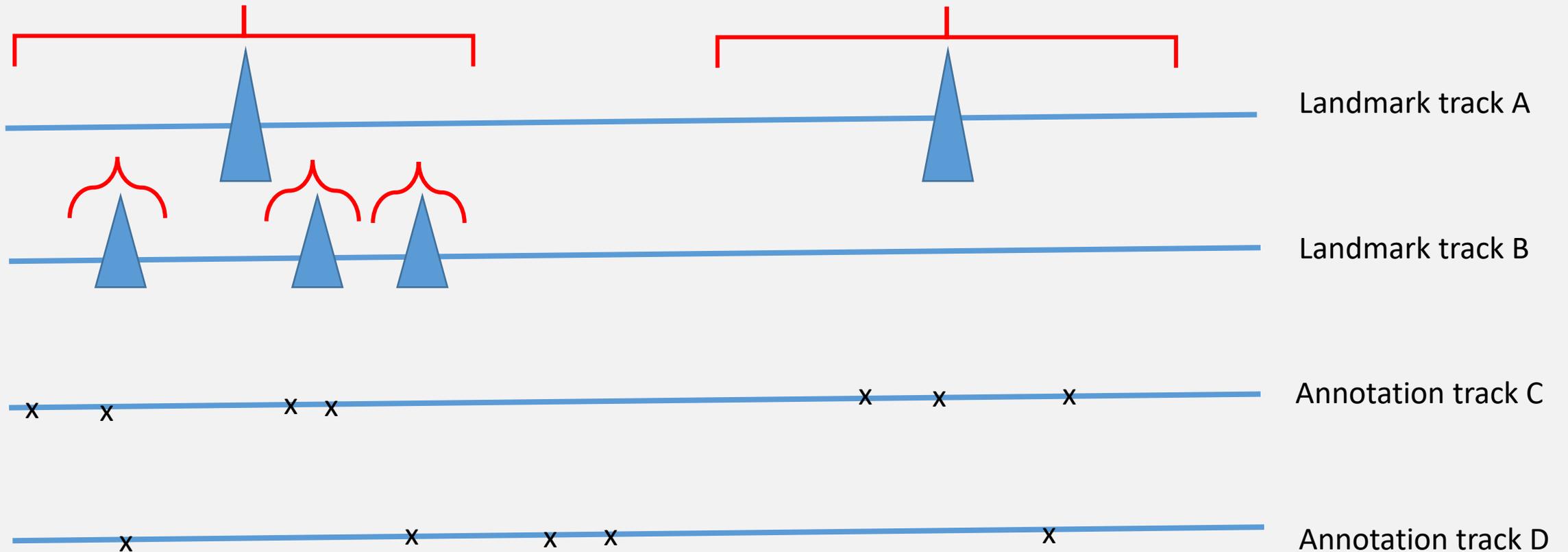
Meta info can be added to tracks. But let's not worry about these here

Equipped with some implicit / automatic normalizations / constraints, e.g. sorted by `#loc`, idempotency and non-overlapping loci on the same track

Landmarks on the same landmark track are non-overlapping, and all annotations can “see” no more than one landmark on the same landmark track

Landmarks can be used for organizing storage & distribution of annotations

Conceptual organization of annotation & landmark tracks



Some operations for the loci type

p before q

p overlap q

p near q ...



p, q : Loc

p can-see r



p : Loc, r : !!t

satisfying “p can-see r whenever (p overlap or near q) & q = r.loc”,

plus maybe other convexity constraints to be thought up

Precise set of operations on loci (e.g. p is-nearer q₁ than q₂) is not important here

But a well-designed set of operations should constraint users from “bad” “expensive” queries, while providing sufficient expressive power for commonly needed queries

Operations on annotations and tracks

$$\frac{e_1 : \{! t_1\}, e_2 : \{! t_2\}}{\cup \{! e_1 \mid x \in e_2\} : \{! t_1\}}$$

$$\frac{}{\{!\} : \{! t\}} \quad \frac{e : !t}{\{! e\} : \{! t\}}$$

$$\frac{e_1, e_2 : \{! t\}}{e_1 \cup e_2 : \{! t\}}$$

Semantics: Same as those for sets, except keep things sorted on #loc & maintains other constraints that we may come up with

Ditto for !t, but maybe ban !!(#loc, #anno)

Plus some set-track conversion ops & syntactic sugars

$$\frac{p : \text{Loc}, e : t}{!(\#loc: p, \#anno: e) : !t}$$

$$\frac{e : !t}{e.\text{loc} : \text{Loc}, e.\text{anno} : t}$$

Let's call this language $\text{NRC}_{\text{genome}}$ in this talk

Common genomic queries in $\text{NRC}_{\text{genome}}$

“Extract from a track R , the annotations in a given region (e.g. 21q22.3) of the genome”

$\{! x \mid x \in R, x.\text{loc overlap } 21\text{q}22.3 \}$

“Extract from the HMMPFAM prediction track, those RBP predictions with $p\text{val} < 1\text{E-}6$ ”

$\{! x \mid x \in \text{HMMPFAM}, x.\text{anno.name} = \text{“RBP”}, x.\text{anno.pval} < 1\text{E-}6 \}$

These queries operate on a single track

They can be executed efficiently, viz. $O(n)$, in $\text{NRC}_{\text{genome}}$

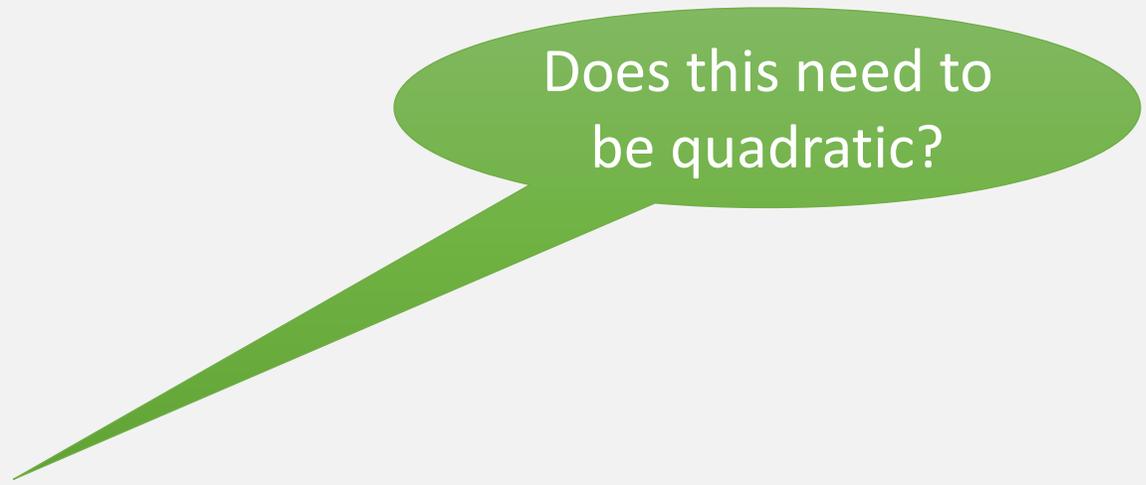
Common genomic queries in $\text{NRC}_{\text{genome}}$

“Extract from the TP53 chip-seq track, those TP53 binding sites with $pval < 1E-6$ and are in promoters of genes”

$\{! x \mid y \in \text{GENES}, x \in \text{TP53},$
 $x.\text{loc}$ before $y.\text{loc},$
 $x.\text{loc}$ near $y.\text{loc},$
 $x.\text{anno.pval} < 1E-6 \}$

This query operates on two tracks

Its “natural” complexity is $O(|\text{GENES}| * |\text{TP53}|)$ in $\text{NRC}_{\text{genome}}$



Does this need to be quadratic?

Common genomic queries in $\text{NRC}_{\text{genome}}$

“Extract from the TP53 and the HDAC1 chip-seq tracks, those TP53 and HDAC1 binding sites that are closest to each other in the promoters of the same genes”

```
{! (#loc: g.loc, #anno: (#name: g.anno.name, #pval: 0, #tp53: u, #hdac1: v))  
| g ∈ GENES,  
  (u, v) ∈ closest { (x, y) | x ∈ TP53, x.loc near g.loc, x.loc before g.loc,  
                           y ∈ HDAC1, y.loc near g.loc, y.loc before g.loc } }
```

This query has complexity $O(|\text{GENES}| * |\text{TP53}| * |\text{HDAC1}|)$ in $\text{NRC}_{\text{genome}}$

Does this need to be cubic?

Common genomic queries in $\text{NRC}_{\text{genome}}$

“Extract from the TP53 and the HDAC1 chip-seq tracks, those TP53 and HDAC1 binding sites that are in the promoters of the same genes”

```
{! u | u ∈ {!  
    (#loc: g.loc, #anno: (#name: g.anno.name, #pval: 0,  
        #tp53: { x | x ∈ TP53, x.loc near g.loc, x.loc before g.loc},  
        #hdac1: { y | y ∈ HDAC1, y.loc near g.loc, y.loc before g.loc} ) )  
    | g ∈ GENES },  
    u.anno.tp53 ≠ {! }, u.anno.hdac1 ≠ {! } }
```

This query has complexity $O(|\text{GENES}| * (|\text{TP53}| + |\text{HDAC1}|))$ in $\text{NRC}_{\text{genome}}$

Does this need to be quadratic?

What is needed? An idea

$$e : \{! t\}, e_1 : \{!! t_1\}, e_2 : \{! t_2\}, \dots, e_k : \{! t_k\}$$
$$\cup \{! e \mid (x_1, x_2, \dots, x_k) \in \in (e_1, e_2, \dots, e_k)\} : \{! t\}$$

Semantics

$$\cup \{! e \mid (x_1, x_2, \dots, x_k) \in \{ (\mathbf{x}_1, \{ \mathbf{x}_2 \mid \mathbf{x}_2 \in \mathbf{e}_2, \mathbf{x}_2.\text{loc can-see } \mathbf{x}_1\}, \dots, \\ \{ \mathbf{x}_k \mid \mathbf{x}_k \in \mathbf{e}_k, \mathbf{x}_k.\text{loc can-see } \mathbf{x}_1\}) \\ \mid \mathbf{x}_1 \in \mathbf{e}_1 \} \}$$

The part in bold is executed for each landmark, considering only annotations which can see that landmark (assuming these are stored with that landmark)

Common genomic queries revisited

“Extract from the TP53 chip-seq track those TP53 binding sites with $pval < 1E-6$ and are in promoters of genes”

$$\{! x \mid y \in \text{GENES}, x \in \text{TP53}, \\ x.\text{loc before } y.\text{loc}, \\ x.\text{loc near } y.\text{loc}, \\ x.\text{anno.pval} < 1E-6 \}$$

GENES is a landmark track

$$\{! x \mid (y, X) \in \in (\text{GENES}, \text{TP53}), x \in X, \\ x.\text{loc before } y.\text{loc}, \\ x.\text{loc near } y.\text{loc}, \\ x.\text{anno.pval} < 1E-6 \}$$

Complexity is maybe $O(|\text{GENES}| * 1\% \text{ of } |\text{TP53}|)$

Common genomic queries revisited

“Extract from the TP53 and the HDAC1 chip-seq tracks, those TP53 and HDAC1 binding sites that are closest to each other in the promoters of the same genes”

```
{! (#loc: g.loc, #anno: (#name: g.anno.name, #pval: 0, #tp53: u, #hdac1: v))  
| g ∈ GENES,  
  (u, v) ∈ closest { (x, y) | x ∈ TP53, x.loc near g.loc, x.loc before g.loc,  
                        y ∈ HDAC1, y.loc near g.loc, y.loc before g.loc } }
```

```
{! (#loc: g.loc, #anno: (#name: g.anno.name, #pval: 0, #tp53: u, #hdac1: v))  
| (g, U, V) ∈∈ (GENES, TP53, HDAC1) ,  
  (u, v) ∈ closest {(x,y) | x ∈ U, x.loc near g.loc, x.loc before g.loc,  
                        y ∈ V, y.loc near g.loc, y.loc before g.loc} }
```

Complexity is maybe $O(|GENES| * (1\% \text{ of } |TP53| * 1\% \text{ of } |HDAC1|))$

Common genomic queries revisited

“Extract from the TP53 and the HDAC1 chip-seq tracks, those TP53 and HDAC1 binding sites that are in the promoters of the same genes”

```
{! (#loc: g.loc, #anno: (#name: g.anno.name, #pval: 0,  
  #tp53: { x | x ∈ TP53, x.loc near g.loc, x.loc before g.loc},  
  #hdac1: { y | y ∈ HDAC1, y.loc near g.loc, y.loc before g.loc } ) )  
| g ∈ GENES}
```

Is it necessary
to process U
and V twice?

```
{! (#loc: g.loc, #anno: (#name: g.anno.name, #pval: 0,  
  #tp53: {! x | x ∈ U, x.loc near g.loc, x.loc before g.loc},  
  #hdac1: {! y | y ∈ V, y.loc near g.loc, y.loc before g.loc}))  
| (g, U, V) ∈ ∈ (GENES, TP53, HDAC1) }
```

Complexity is maybe $O(|GENES| * (1\% \text{ of } |TP53| + 1\% \text{ of } |HDAC1|))$

A better idea?

$e : \{! t\}, e_1 : \{!! t_1\}, e_2 : \{! t_2\}, \dots, e_k : \{! t_k\}, \gamma_1 : \text{bool}, \dots, \gamma_k : \text{bool}$

$\cup \{! e \mid x_1 \in e_1 \text{ st } \gamma_1, X_2 \subseteq e_2 \ni x_2 \text{ st } \gamma_2, \dots, X_k \subseteq e_k \ni x_k \text{ st } \gamma_k\} : \{! t\}$

$FV(\gamma_j) \setminus \{x_1, x_j\} \subseteq FV(\cup \{! e \mid x_1 \in e_1, X_2 \subseteq e_2 \ni x_2 \text{ st } \gamma_2, \dots, X_k \subseteq e_k \ni x_k \text{ st } \gamma_k\})$, and $FV(e) \cap \{x_2, \dots, x_k\} = \{\}$

Semantics

$\cup \{! e \mid (x_1, X_2, \dots, X_k) \in \{ (\mathbf{x}_1, \{ \mathbf{x}_2 \mid \mathbf{x}_2 \in e_2, \mathbf{x}_2.\text{loc can-see } \mathbf{x}_1, \gamma_2 \}), \dots, \{ \mathbf{x}_k \mid \mathbf{x}_k \in e_k, \mathbf{x}_k.\text{loc can-see } \mathbf{x}_1, \gamma_k \}) \mid \mathbf{x}_1 \in e_1, \gamma_1 \} \}$

The part in bold is executed for each landmark, considering only annotations which can see that landmark (assuming these are stored with that landmark)

Common genomic queries revisited again

“Extract from the TP53 and the HDAC1 chip-seq tracks, those TP53 and HDAC1 binding sites that are in the promoters of the same genes”

```
{! (#loc: g.loc, #anno: (#name: g.anno.name, #pval: 0,  
  #tp53: { x | x ∈ TP53, x.loc near g.loc, x.loc before g.loc},  
  #hdac1: { y | y ∈ HDAC1, y.loc near g.loc, y.loc before g.loc } ) )  
| g ∈ GENES}
```

```
{! (#loc: g.loc, #anno: (#name: g.anno.name, #pval: 0, #tp53: U, #hdac1: V))  
| g ∈ GENES,  
  U ⊆ TP53 ∃ u st u.loc near g.loc & u.loc before g.loc ,  
  V ⊆ HDAC1 ∃ v st v.loc near g.loc, v.loc before g.loc }
```

Complexity is maybe $O(|GENES| * (1\% \text{ of } |TP53| + 1\% \text{ of } |HDAC1|))$

Common genomic queries revisited again

“Extract from the TP53 and the HDAC1 chip-seq tracks, those TP53 and HDAC1 binding sites that are closest to each other in the promoters of the same genes”

```
{! (#loc: g.loc, #anno: (#name: g.anno.name, #pval: 0, #tp53: u, #hdac1: v))  
| g ∈ GENES,  
  (u, v) ∈ closest { (x, y) | x ∈ TP53, x.loc near g.loc, x.loc before g.loc,  
                        y ∈ HDAC1, y.loc near g.loc, y.loc before g.loc } }
```

```
{! (#loc: g.loc, #anno: (#name: g.anno.name, #pval: 0, #tp53: x, #hdac1: y))  
| g ∈ GENES, U ⊆ TP53 ∃ u st u.loc near g.loc & u.loc before g.loc,  
  V ⊆ HDAC1 ∃ v st v.loc near g.loc & v.loc before g.loc ,  
  (x, y) ∈ closest {(u,v) | u ∈ U, v ∈ V} }
```

Complexity is maybe $O(|GENES| * (1\% \text{ of } |TP53| * 1\% \text{ of } |HDAC1|))$

And this idea? It is really a syntactic sugar

$$e : \{! t\}, e_1 : \{!! t_1\}, e_2 : \{! t_2\}, \dots, e_k : \{! t_k\}, \gamma_1 : \text{bool}, \dots, \gamma_k : \text{bool}$$
$$\cup \{! e \mid x_1 \in e_1 \text{ st } \gamma_1, x_2 \in e_2 \text{ st } \gamma_2, \dots, x_k \in e_k \text{ st } \gamma_k\} : \{! t\}$$
$$FV(\gamma_j) \setminus \{x_1, x_j\} \subseteq FV(\cup \{! e \mid x_1 \in e_1, x_2 \in e_2 \text{ st } \gamma_2, \dots, x_k \in e_k \text{ st } \gamma_k\})$$

Semantics

$$\cup \{! e \mid (x_1, X_2, \dots, X_k) \in \{ (\mathbf{x}_1, \{ \mathbf{x}_2 \mid \mathbf{x}_2 \in e_2, \mathbf{x}_2.\text{loc can-see } x_1, \gamma_2 \}, \dots, \\ \{ \mathbf{x}_k \mid \mathbf{x}_k \in e_k, \mathbf{x}_k.\text{loc can-see } x_1, \gamma_k \}) \mid x_1 \in e_1, \gamma_1 \}, \\ x_2 \in X_2, \dots, x_k \in X_k \}$$

The part in bold is executed for each landmark, considering only annotations which can see that landmark (assuming these are stored with that landmark)

Common genomic queries revisited again

“Extract from the TP53 chip-seq track those TP53 binding sites with pval < 1E-6 and are in promoters of genes”

```
{! x | y ∈ GENES, x ∈ TP53,  
      x.loc before y.loc,  
      x.loc near y.loc,  
      x.anno.pval < 1E-6 }
```

```
{! x | y ∈ GENES,  
      x ∈ TP53 st x.loc before y.loc &  
                x.loc near y.loc & x.anno.pval < 1E-6 }
```

Complexity is maybe $O(|GENES| * 1\% \text{ of } |TP53|)$

GENES is a landmark track

Implementing “synchronized” processing of multiple lists / tracks

$\text{lzip}: (t_1 \rightarrow \text{bool}) * (t_1 * t_2 \rightarrow \text{bool}) * (t_1 * t_2 \rightarrow \text{bool}) * (t_2 * t' \rightarrow t') * (t_1 * t' \rightarrow t') * (t' \rightarrow \{t\}) * t' * t'$
 $\rightarrow \{t_1\} * \{t_2\} \rightarrow \{t\}$

$\text{lzip} (\text{sx}, \text{sy}, \text{ay}, \text{h}, \text{g}, \text{f}, \text{a}, \text{e}) (\{\}, Y) = \text{f a}$

$\text{lzip} (\text{sx}, \text{sy}, \text{ay}, \text{h}, \text{g}, \text{f}, \text{a}, \text{e}) (X, \{\}) = \text{f a}$

$\text{lzip} (\text{sx}, \text{sy}, \text{ay}, \text{h}, \text{g}, \text{f}, \text{a}, \text{e}) (x::X, y::Y) =$

if $\text{sx}(x)$

then if $\text{sy}(x, y)$

then if $\text{ay}(x, y)$

then $\text{lzip} (\text{sx}, \text{sy}, \text{ay}, \text{h}, \text{g}, \text{f}, \text{h}(y, \text{g}(x, \text{a})), \text{e}) (x::X, Y)$

else $\text{lzip} (\text{sx}, \text{sy}, \text{ay}, \text{h}, \text{g}, \text{f}, \text{g}(x, \text{a}), \text{e}) (x::X, Y)$

else $\text{f} (\text{g}(x, \text{a})) \cup \text{lzip} (\text{sx}, \text{sy}, \text{ay}, \text{h}, \text{g}, \text{f}, \text{e}, \text{e}) (X, y::Y)$

else $\text{f a} \cup \text{lzip} (\text{sx}, \text{sy}, \text{ay}, \text{h}, \text{g}, \text{f}, \text{e}, \text{e}) (X, y::Y)$

At every step, either x or y gets shifted. So complexity is $O(|X| + |Y| * \alpha)$, where α is complexity of sx , sy , etc.

Implementing $\cup\{! e \mid x_1 \in e_1 \text{ st } \gamma_1, X_2 \subseteq e_2 \ni x_2 \text{ st } \gamma_2, \dots, X_k \subseteq e_k \ni x_k \text{ st } \gamma_k\}$

$\cup\{! e \mid x_1 \in e_1 \text{ st } \gamma_1, X_2 \subseteq e_2 \ni x_2 \text{ st } \gamma_2\} :=$

$\text{lzip}(sx, sy, ay, h, g, f, (\{\!\!\}, \{\!\}), (\{\!\!\}, \{\!\})) (e_1, e_2)$ where

$sx(x_1) := \gamma_1,$

$ay(x_1, x_2) := x_2.\text{loc can-see } x_1 \text{ \& } \gamma_2,$

$sy(x_1, x_2) := x_2.\text{loc before } x_1.\text{loc or } ay(x_1, x_2),$

$h(x_2, (X_1, X_2)) := (X_1, X_2 \cup \{\!\ x_2\ \}),$

$g(x_1, (X_1, X_2)) := (X_1 \cup \{\!\!\ x_1\ \}, X_2),$

$f(X_1, X_2) := \cup\{! e \mid x_1 \in X_1\};$

Synchronized scan

```
lzip: (t1 → bool) * (t1 * t2 → bool) * (t1 * t2 → bool) * (t2 * t' → t') * (t1 * t' → t') * (t' → {t}) * t * t'
      → {t1} * {t2} → {t}
```

```
lzip (sx, sy, ay, h, g, f, a, e) ({}, Y) = f a
```

```
lzip (sx, sy, ay, h, g, f, a, e) (X, {}) = f a
```

```
lzip (sx, sy, ay, h, g, f, a, e) (x::X, y::Y) =
```

```
  if sx(x)
```

```
  then if sy(x, y)
```

```
    then if ay(x, y)
```

```
      then lzip (sx, sy, ay, h, g, f, h(y, g(x, a)), e) (x::X, Y)
```

```
      else lzip (sx, sy, ay, h, g, f, g(x, a), e) (x::X, Y)
```

```
    else f (g(x, a)) ∪ lzip (sx, sy, ay, h, g, f, e, e) (X, y::Y)
```

```
  else f a ∪ lzip (sx, sy, ay, h, g, f, e, e) (X, y::Y)
```

A nice property of $\cup\{! e \mid x_1 \in e_1 \text{ st } \gamma_1, X_2 \subseteq e_2 \ni x_2 \text{ st } \gamma_2, \dots, X_k \subseteq e_k \ni x_k \text{ st } \gamma_k\}$

$\cup\{! e \mid x_1 \in e_1 \text{ st } \gamma_1, X_2 \subseteq e_2 \ni x_2 \text{ st } \gamma_2\} :=$

$\text{lzip}(\dots)(e_1, e_2)$ where ...

is a homomorphism on e_1 . Thus

$$\begin{aligned} & \cup\{! e \mid x_1 \in e_1 \in \{!! o_1, \dots, o_k\} \text{ st } \gamma_1, X_2 \subseteq e_2 \ni x_2 \text{ st } \gamma_2\} \\ = & \cup\{! e \mid x_1 \in e_1 \in \{!! o_1\} \text{ st } \gamma_1, X_2 \subseteq e_2 \ni x_2 \text{ st } \gamma_2\} \cup \dots \cup \\ & \cup\{! e \mid x_1 \in e_1 \in \{!! o_k\} \text{ st } \gamma_1, X_2 \subseteq e_2 \ni x_2 \text{ st } \gamma_2\} \end{aligned}$$

```
lzip: (t1 → bool) * (t1 * t2 → bool) * (t1 * t2 → bool) * (t2 * t' → t') * (t1 * t' → t') * (t' → {t}) * t * t'
      → {t1} * {t2} → {t}
```

```
lzip (sx, sy, ay, h, g, f, a, e) ({}, Y) = f a
```

```
lzip (sx, sy, ay, h, g, f, a, e) (X, {}) = f a
```

```
lzip (sx, sy, ay, h, g, f, a, e) (x::X, y::Y) =
```

```
  if sx(x)
```

```
  then if sy(x, y)
```

```
    then if ay(x,y)
```

```
      then lzip (sx, sy, ay, h, g, f, h(y, g(x, a)), e) (x::X, Y)
```

```
      else lzip (sx, sy, ay, h, g, f, g(x, a), e) (x::X, Y)
```

```
    else f (g(x, a)) ∪ lzip (sx, sy, ay, h, g, f, e, e) (X, y::Y)
```

```
  else f a ∪ lzip (sx, sy, ay, h, g, f, e, e) (X, y::Y)
```

When annotations on track e_2 are “clustered” (i.e. stored with) the specific landmarks on track e_1 these annotations “can see”, each $\cup\{! e \mid x_1 \in e_1 \in \{!! o_j\} \text{ st } \gamma_1, X_2 \subseteq e_2 \ni x_2 \text{ st } \gamma_2\}$ can be run in parallel on each cluster

Some optimization rules

$$\cup \{! \text{ if } \varphi \text{ then } e \text{ else } \{! \} \mid x_1 \in e_1 \text{ st } \gamma_1, x_2 \in e_2 \text{ st } \gamma_2, \dots, x_k \in e_k \text{ st } \gamma_k \}$$
$$\Downarrow$$
$$\Downarrow x_j \in \text{FV}(\varphi) \text{ and } \text{FV}(\varphi) \cap \{x_1, \dots, x_k\} \subseteq \{x_1, x_j\}$$
$$\Downarrow$$
$$\cup \{! e \mid x_1 \in e_1 \text{ st } \gamma_1, x_2 \in e_2 \text{ st } \gamma_2, \dots, x_j \in e_j \text{ st } \gamma_j \ \& \ \varphi, \dots, x_k \in e_k \text{ st } \gamma_k \}$$

And φ is a “positive” condition on loci in both rules

$$\cup \{! \cup \{! \text{ if } \varphi \text{ then } e \text{ else } \{! \} \mid x \in X_j \}$$
$$\mid x_1 \in e_1 \text{ st } \gamma_1, X_2 \subseteq e_2 \ni x_2 \text{ st } \gamma_2, \dots, X_k \subseteq e_k \ni x_2 \text{ st } \gamma_k \}$$
$$\Downarrow$$
$$\Downarrow x \in \text{FV}(\varphi) \text{ and } \text{FV}(\varphi) \cap \{x, x_1, X_2, \dots, X_k\} \subseteq \{x_1, x\}$$
$$\Downarrow$$
$$\cup \{! \cup \{! e \mid x \in X_j \}$$
$$\mid x_1 \in e_1 \text{ st } \gamma_1,$$
$$X_2 \subseteq e_2 \ni x_2 \text{ st } \gamma_2, \dots, X_j \subseteq e_j \ni x_j \text{ st } \gamma_j \ \& \ \varphi[x_j/x], \dots, X_k \subseteq e_k \ni x_2 \text{ st } \gamma_k \}$$

And ...

$\cup\{! \cup\{! \text{ if } \varphi \text{ then } e \text{ else } \{! \} \mid x_2 \in e_2 \} \mid x_1 \in e_1\}$

\Downarrow

$\Downarrow e_1: !t_1 \ \& \ e_2: !t_2 \ \&$

$\Downarrow x_1 \notin \text{FV}(e_2) \ \& \ \text{FV}(\varphi) \cap \{x_1, x_2\} = \{x_1, x_2\} \ \&$

$\Downarrow \varphi \text{ is a positive condition on loci of } x_1, x_2$

\Downarrow

$\cup\{! \cup\{! \text{ if } \varphi \text{ then } e \text{ else } \{! \} \mid x_1 \in e_1 \text{ st true}, x_2 \in e_2 \text{ st true} \}$

So a user does not need to worry about when to use $\cup\{! e \mid x_1 \in e_1, \dots\}$

In fact, ...

It is not necessary for a user to use $\{!$, $\{!!$, etc.

These can be inferred by a simple type system

And transformed into synchronized/parallel scans by an optimizer

