

# Analysis of the Subfunctionalization Model for the Fate of Gene Duplicates

Tristan L. Stark <sup>1</sup>  
Malgorzata O'Reilly <sup>1</sup>  
Barbara Holland <sup>1</sup>  
David Liberles <sup>2</sup>

<sup>1</sup>University of Tasmania

<sup>2</sup>Temple University, Philadelphia

November 18, 2015



We consider the evolution of a pair of gene duplicates, each with  $z$  regulatory regions and a coding region.



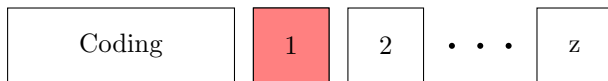
We consider the evolution of a pair of gene duplicates, each with  $z$  regulatory regions and a coding region.

Under the subfunctionalization model a null mutation can fix

We consider the evolution of a pair of gene duplicates, each with  $z$  regulatory regions and a coding region.

Under the subfunctionalization model a null mutation can fix

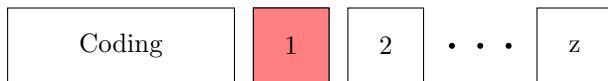
- in any of the  $z$  regulatory regions of either copy. We assume this occurs at equal Poisson rate  $u_r$  for all  $2z$  regions.



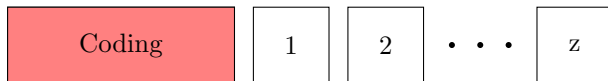
We consider the evolution of a pair of gene duplicates, each with  $z$  regulatory regions and a coding region.

Under the subfunctionalization model a null mutation can fix

- in any of the  $z$  regulatory regions of either copy. We assume this occurs at equal Poisson rate  $u_r$  for all  $2z$  regions.



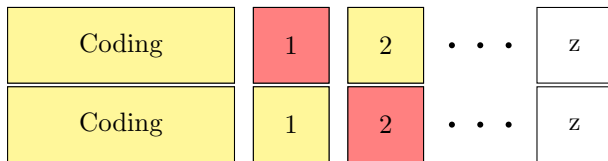
- in the coding region of either gene. We assume this occurs at Poisson rate  $u_c$  for each gene.



As mutations build up in the two copies, one of two possible fates eventually occurs

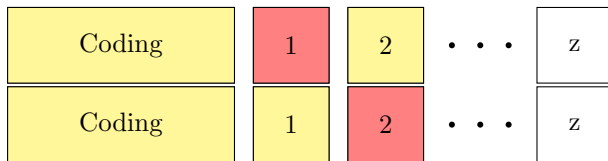
As mutations build up in the two copies, one of two possible fates eventually occurs

- Subfunctionalization

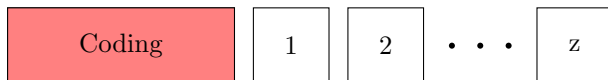


As mutations build up in the two copies, one of two possible fates eventually occurs

- Subfunctionalization



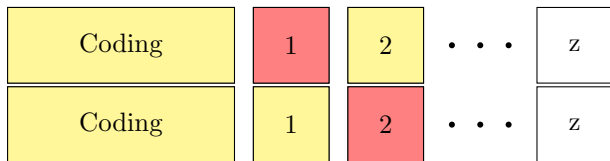
- Pseudogenization, or gene loss



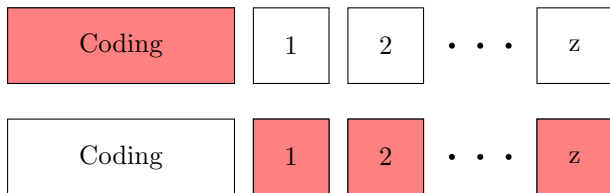


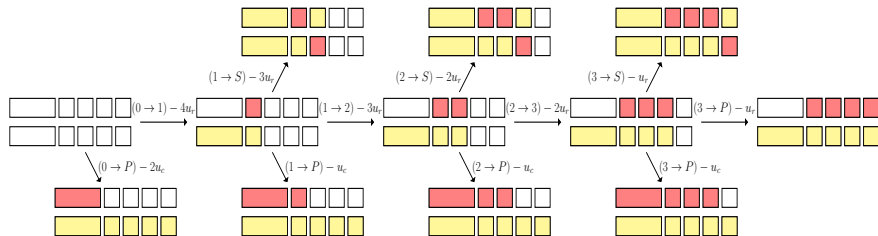
As mutations build up in the two copies, one of two possible fates eventually occurs

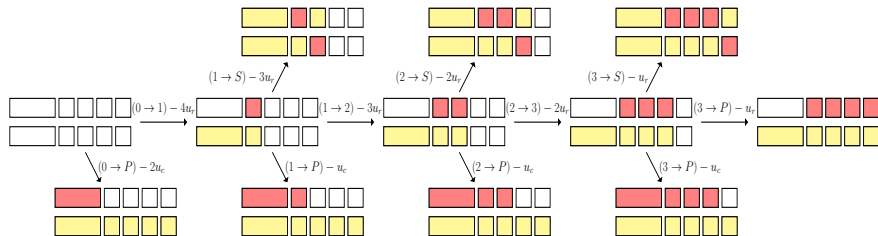
- Subfunctionalization



- Pseudogenization, or gene loss

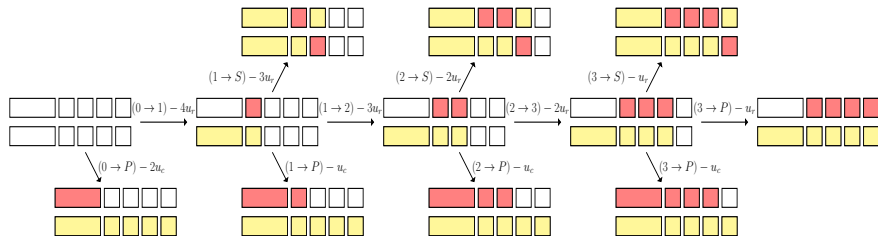






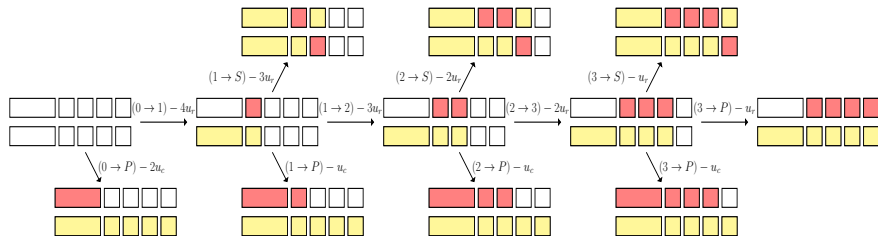
## Notice

- Initial rate of Pseudogenization is  $2u_c$



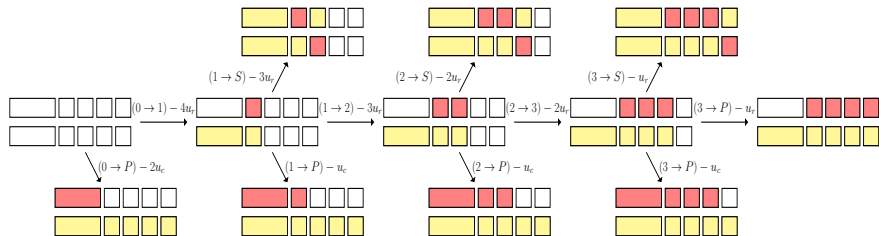
## Notice

- Initial rate of Pseudogenization is  $2u_c$
- After first mutation, drops to  $u_c$



## Notice

- Initial rate of Pseudogenization is  $2u_c$
- After first mutation, drops to  $u_c$
- At final mutation, increases to  $u_c + u_r$ .



## Notice

- Initial rate of Pseudogenization is  $2u_c$
- After first mutation, drops to  $u_c$
- At final mutation, increases to  $u_c + u_r$ .
- Rate of Subfunctionalization equals rate of transition to  $i + 1$  equals  $(z - i)u_r$ .



Model is given by  $\mathbf{Q} = [q_{ij}]$  where



Model is given by  $\mathbf{Q} = [q_{ij}]$  where

$$q_{ij} = \begin{cases} 2u_c & \text{if } i = 0, j = P \\ 2zu_r & \text{if } i = 0, j = 1 \\ u_c & \text{if } 1 \leq i \leq z - 2, j = P \\ (z - i)u_r & \text{if } 1 \leq i \leq z - 2, j = i + 1 \text{ or } j = S \\ u_r + u_c & \text{if } i = z - 1, j = P \\ u_r & \text{if } i = z - 1, j = S. \end{cases} \quad (1)$$





Model is given by  $\mathbf{Q} = [q_{ij}]$  where

$$q_{ij} = \begin{cases} 2u_c & \text{if } i = 0, j = P \\ 2zu_r & \text{if } i = 0, j = 1 \\ u_c & \text{if } 1 \leq i \leq z - 2, j = P \\ (z - i)u_r & \text{if } 1 \leq i \leq z - 2, j = i + 1 \text{ or } j = S \\ u_r + u_c & \text{if } i = z - 1, j = P \\ u_r & \text{if } i = z - 1, j = S. \end{cases} \quad (1)$$

The structure of this CTMC is much like those that give rise to the phase-type distribution.



For CTMCs of this structure, it is convenient to write

$$\mathbf{Q} = \left[ \begin{array}{c|c} \mathbf{Q}^* & \mathbf{V} \\ \hline \mathbf{0} & \mathbf{0} \end{array} \right], \quad (2)$$

where  $\mathbf{Q}^*$  contains the entries corresponding to transitions between transient states, and  $\mathbf{V}$  transitions to absorbing states.

Its possible to exploit the phase-type-like structure of our chain to derive many measures of interest.

Its possible to exploit the phase-type-like structure of our chain to derive many measures of interest.

### Probability density of absorption

$$f_i(t) = \underline{\mathbf{e}}_i e^{\mathbf{Q}^* t} \mathbf{V} \mathbf{1} \quad (3)$$

Its possible to exploit the phase-type-like structure of our chain to derive many measures of interest.

## Probability density of absorption

$$f_i(t) = \underline{\mathbf{e}}_i e^{\mathbf{Q}^* t} \mathbf{V} \mathbf{1} \quad (3)$$

## Cumulative distribution function

$$\begin{aligned} F_i(t) &= \int_{u=0}^t f_i(u) du \\ &= \int_0^t \underline{\mathbf{e}}_i e^{\mathbf{Q}^* u} \mathbf{V} \mathbf{1} du \end{aligned}$$



Its possible to exploit the phase-type-like structure of our chain to derive many measures of interest.

### Probability density of absorption

$$f_i(t) = \underline{\mathbf{e}}_i e^{\mathbf{Q}^* t} \mathbf{V} \mathbf{1} \quad (3)$$

### Cumulative distribution function

$$\begin{aligned} F_i(t) &= \int_{u=0}^t f_i(u) du \\ &= \int_0^t \underline{\mathbf{e}}_i e^{\mathbf{Q}^* u} \mathbf{V} \mathbf{1} du \end{aligned}$$

Using the fact that  $\mathbf{Q}^* \mathbf{1} + \mathbf{V} \mathbf{1} = 0$  its easy to show



Its possible to exploit the phase-type-like structure of our chain to derive many measures of interest.

### Probability density of absorption

$$f_i(t) = \underline{\mathbf{e}}_i e^{\mathbf{Q}^* t} \mathbf{V} \mathbf{1} \quad (3)$$

### Cumulative distribution function

$$\begin{aligned} F_i(t) &= \int_{u=0}^t f_i(u) du \\ &= \int_0^t \underline{\mathbf{e}}_i e^{\mathbf{Q}^* u} \mathbf{V} \mathbf{1} du \end{aligned}$$

Using the fact that  $\mathbf{Q}^* \mathbf{1} + \mathbf{V} \mathbf{1} = 0$  its easy to show

$$F_i(t) = 1 - \underline{\mathbf{e}}_i e^{\mathbf{Q}^* t} \mathbf{1}. \quad (4)$$



With the density and cumulative distribution functions, we're able to derive results for various measures

- Probability of absorption into  $j \in \{S, P\}$

$$\begin{aligned} p_{i,j} &= \int_{t=0}^{\infty} \underline{\mathbf{e}}_i e^{\mathbf{Q}^* t} \mathbf{V}_j dt \\ &= -\underline{\mathbf{e}}_i (\mathbf{Q}^*)^{(-1)} \mathbf{V}_j \end{aligned} \quad (5)$$

- The  $k^{\text{th}}$  moment of time until absorption

$$\begin{aligned} m_i^{(k)} &= \int_{t=0}^{\infty} t^k \underline{\mathbf{e}}_i e^{\mathbf{Q}^* t} \mathbf{V} \mathbf{1} dt \\ &= (-1)^k k! \underline{\mathbf{e}}_i (\mathbf{Q}^*)^{(-k)} \mathbf{1}, \end{aligned} \quad (6)$$

- Variance of time until absorption

$$\text{var}_i = m_i^{(2)} - (m_i)^2. \quad (7)$$





When there are several absorbing states, often interested in the cause-specific hazard rate

When there are several absorbing states, often interested in the cause-specific hazard rate

## Cause-specific hazard rate

$$\begin{aligned}\lambda_{ij}(t) &= \lim_{h \rightarrow 0^+} \frac{P(t < T_{\{S,P\}} < t + h, X(T_{\{S,P\}}) = j | T_{\{S,P\}} > t, X(0) = i)}{h} \\ &= \frac{f_{ij}(t)}{1 - F_i(t)} = \frac{\mathbf{e}_i \mathbf{e}^{\mathbf{Q}^* t} \mathbf{V}_j}{\mathbf{e}_0 \mathbf{e}^{\mathbf{Q}^* t} \mathbf{1}}\end{aligned}\quad (8)$$

$$f_i(t) = \sum_{j \in \{S,P\}} f_{ij}(t),$$

$$\lambda_i(t) = \sum_{j \in \{S,P\}} \lambda_{ij}(t).\quad (9)$$



The numerical states in our chain are essentially a pure-birth process (recall we track the number of mutations to have occurred).

The numerical states in our chain are essentially a pure-birth process (recall we track the number of mutations to have occurred).

Think of the process eventually reaching the state  $z - 1$ , with no possibility of transition from  $i$  to  $j < i$ .

The numerical states in our chain are essentially a pure-birth process (recall we track the number of mutations to have occurred).

Think of the process eventually reaching the state  $z - 1$ , with no possibility of transition from  $i$  to  $j < i$ .

Recall

- $q_{z-1,P} = u_c + u_r,$

The numerical states in our chain are essentially a pure-birth process (recall we track the number of mutations to have occurred).

Think of the process eventually reaching the state  $z - 1$ , with no possibility of transition from  $i$  to  $j < i$ .

Recall

- $q_{z-1,P} = u_c + u_r,$
- $q_{z-1,S} = u_r.$

The numerical states in our chain are essentially a pure-birth process (recall we track the number of mutations to have occurred).

Think of the process eventually reaching the state  $z - 1$ , with no possibility of transition from  $i$  to  $j < i$ .

Recall

- $q_{z-1,P} = u_c + u_r$ ,
- $q_{z-1,S} = u_r$ .

Since hazard rate assumes process is not absorbed, as  $t$  becomes large  $X(t)$  is almost certainly  $z - 1$ .

The numerical states in our chain are essentially a pure-birth process (recall we track the number of mutations to have occurred).

Think of the process eventually reaching the state  $z - 1$ , with no possibility of transition from  $i$  to  $j < i$ .

Recall

- $q_{z-1,P} = u_c + u_r$ ,
- $q_{z-1,S} = u_r$ .

Since hazard rate assumes process is not absorbed, as  $t$  becomes large  $X(t)$  is almost certainly  $z - 1$ .

So  $\lim_{t \rightarrow \infty} \lambda_{ij}$  is surely  $u_c + u_r$  for  $j = P$   $u_r$  for  $j = S$ .



We define the following:

## Some events

- The event that processes has not been absorbed by time  $t$ , but is absorbed by later time  $t + h$

$$A_t^h = \{t < T_{\{S,P\}} < t + h\}.$$

We define the following:

## Some events

- The event that processes has not been absorbed by time  $t$ , but is absorbed by later time  $t + h$

$$A_t^h = \{t < T_{\{S,P\}} < t + h\}.$$

- The event that process has not been absorbed by time  $t$

$$B_t = \{T_{\{S,P\}} > t\}.$$

We define the following:

## Some events

- The event that process has not been absorbed by time  $t$ , but is absorbed by later time  $t + h$

$$A_t^h = \{t < T_{\{S,P\}} < t + h\}.$$

- The event that process has not been absorbed by time  $t$

$$B_t = \{T_{\{S,P\}} > t\}.$$

- The event that the process has entered state  $z - 1$  by time  $t$

$$C_t = \{T_{z-1} \leq t\}.$$

We define the following:

## Some events

- The event that process has not been absorbed by time  $t$ , but is absorbed by later time  $t + h$

$$A_t^h = \{t < T_{\{S,P\}} < t + h\}.$$

- The event that process has not been absorbed by time  $t$

$$B_t = \{T_{\{S,P\}} > t\}.$$

- The event that the process has entered state  $z - 1$  by time  $t$

$$C_t = \{T_{z-1} \leq t\}.$$

- The event that the process has not entered state  $z - 1$  by time  $t$

$$\bar{C}_t = \{T_{z-1} > t\}$$



First, note that

$$\lim_{t \rightarrow \infty} P(C_t | B_t, X(0) = i) = 1 \quad \text{and} \quad \lim_{t \rightarrow \infty} P(\bar{C}_t | B_t, X(0) = i) = 0,$$

First, note that

$$\lim_{t \rightarrow \infty} P(C_t | B_t, X(0) = i) = 1 \quad \text{and} \quad \lim_{t \rightarrow \infty} P(\bar{C}_t | B_t, X(0) = i) = 0,$$

Now,

$$\lim_{t \rightarrow \infty} \lambda_{ij}(t) = \lim_{t \rightarrow \infty} \lim_{h \rightarrow 0^+} \frac{P(A_t^h, X(T_{\{S,P\}}) = j | B_t, X(0) = i)}{h}$$



First, note that

$$\lim_{t \rightarrow \infty} P(C_t | B_t, X(0) = i) = 1 \quad \text{and} \quad \lim_{t \rightarrow \infty} P(\bar{C}_t | B_t, X(0) = i) = 0,$$

Now,

$$\lim_{t \rightarrow \infty} \lambda_{ij}(t) = \lim_{t \rightarrow \infty} \lim_{h \rightarrow 0^+} \frac{P(A_t^h, X(T_{\{S,P\}}) = j | B_t, X(0) = i)}{h}$$

by law of total probability and two limits above this becomes

$$\lim_{t \rightarrow \infty} \lim_{h \rightarrow 0^+} \frac{P(A_t^h, X(T_{\{S,P\}}) = j | B_t, X(0) = i, C_t)}{h}$$



First, note that

$$\lim_{t \rightarrow \infty} P(C_t | B_t, X(0) = i) = 1 \quad \text{and} \quad \lim_{t \rightarrow \infty} P(\bar{C}_t | B_t, X(0) = i) = 0,$$

Now,

$$\lim_{t \rightarrow \infty} \lambda_{ij}(t) = \lim_{t \rightarrow \infty} \lim_{h \rightarrow 0^+} \frac{P(A_t^h, X(T_{\{S,P\}}) = j | B_t, X(0) = i)}{h}$$

by law of total probability and two limits above this becomes

$$\lim_{t \rightarrow \infty} \lim_{h \rightarrow 0^+} \frac{P(A_t^h, X(T_{\{S,P\}}) = j | B_t, X(0) = i, C_t)}{h}$$

Which simplifies to

$$\lim_{t \rightarrow \infty} \lim_{h \rightarrow 0^+} \frac{P(X(t+h) = j | X(t) = z-1)}{h}$$

(10)



First, note that

$$\lim_{t \rightarrow \infty} P(C_t | B_t, X(0) = i) = 1 \quad \text{and} \quad \lim_{t \rightarrow \infty} P(\bar{C}_t | B_t, X(0) = i) = 0,$$

Now,

$$\lim_{t \rightarrow \infty} \lambda_{ij}(t) = \lim_{t \rightarrow \infty} \lim_{h \rightarrow 0^+} \frac{P(A_t^h, X(T_{\{S,P\}}) = j | B_t, X(0) = i)}{h}$$

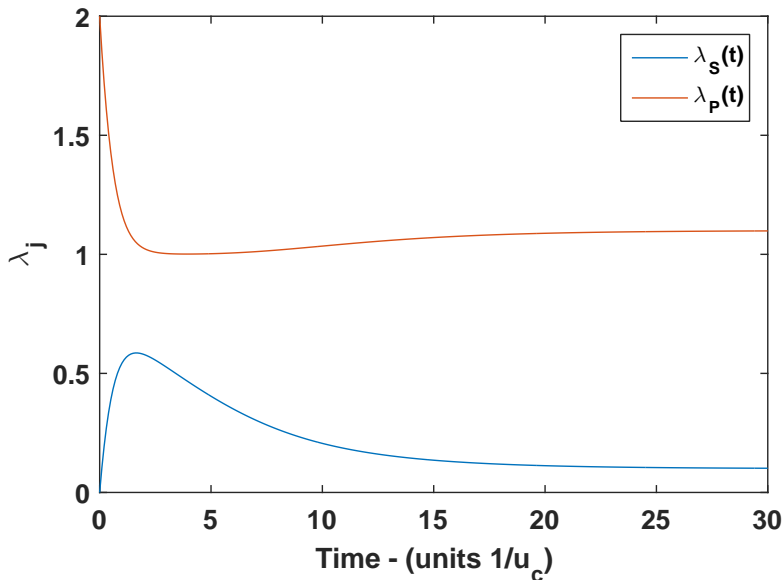
by law of total probability and two limits above this becomes

$$\lim_{t \rightarrow \infty} \lim_{h \rightarrow 0^+} \frac{P(A_t^h, X(T_{\{S,P\}}) = j | B_t, X(0) = i, C_t)}{h}$$

Which simplifies to

$$\lim_{t \rightarrow \infty} \lim_{h \rightarrow 0^+} \frac{P(X(t+h) = j | X(t) = z-1)}{h} \tag{10}$$

By Markov Property we can drop the  $t$ 's, and we're left with  $q_{z-1,j}$





We might be interested in the rate of absorption into state  $P$  at time  $t$  conditional only on not having been absorbed into  $P$ .



We might be interested in the rate of absorption into state  $P$  at time  $t$  conditional only on not having been absorbed into  $P$ .

We define the following rate

## Pseudogenization rate

$$\begin{aligned}
 h_P^z(t) &= \lim_{h \rightarrow 0^+} \frac{P(t < T_P < t + h | T_P > t, X(0) = 0)}{h} \\
 &= \frac{f(t)}{1 - F(t)} \\
 &= \frac{\underline{\mathbf{e}}_0 e^{\mathbf{Q}^* t} \mathbf{V}_P}{1 - \int_{u=0}^t \underline{\mathbf{e}}_0 e^{\mathbf{Q}^* u} \mathbf{V}_P du} \\
 &= \frac{\underline{\mathbf{e}}_0 e^{\mathbf{Q}^* t} \underline{\mathbf{V}}_P}{1 - \underline{\mathbf{e}}_0 (e^{\mathbf{Q}^* t} - \mathbf{I}) (\mathbf{Q}^*)^{(-1)} \mathbf{V}_P}.
 \end{aligned} \tag{11}$$

Here  $T_P$  is RV tracking time to absorption into  $P$ , and could be infinity.



Intuitively, expect pseudogenization rate to go to 0 as  $t \rightarrow \infty$ .



Intuitively, expect pseudogenization rate to go to 0 as  $t \rightarrow \infty$ .  
This is easily proved using fact that

$$\lim_{t \rightarrow \infty} e^{\mathbf{Q}^* t} = 0, \quad (12)$$



Intuitively, expect pseudogenization rate to go to 0 as  $t \rightarrow \infty$ .  
This is easily proved using fact that

$$\lim_{t \rightarrow \infty} e^{\mathbf{Q}^* t} = 0, \quad (12)$$

Notice

$$\lim_{t \rightarrow \infty} h_P(t) = \frac{\lim_{t \rightarrow \infty} \underline{\mathbf{e}}_0 e^{\mathbf{Q}^* t} \underline{\mathbf{v}}_P}{1 - \lim_{t \rightarrow \infty} \underline{\mathbf{e}}_0 (e^{\mathbf{Q}^* t} - \mathbf{I}) (\mathbf{Q}^*)^{(-1)} \underline{\mathbf{v}}_P}$$



Intuitively, expect pseudogenization rate to go to 0 as  $t \rightarrow \infty$ .  
This is easily proved using fact that

$$\lim_{t \rightarrow \infty} e^{\mathbf{Q}^* t} = 0, \quad (12)$$

Notice

$$\begin{aligned} \lim_{t \rightarrow \infty} h_P(t) &= \frac{\lim_{t \rightarrow \infty} \underline{\mathbf{e}}_0 e^{\mathbf{Q}^* t} \underline{\mathbf{v}}_P}{1 - \lim_{t \rightarrow \infty} \underline{\mathbf{e}}_0 (e^{\mathbf{Q}^* t} - \mathbf{I}) (\mathbf{Q}^*)^{(-1)} \underline{\mathbf{v}}_P} \\ &= \frac{0}{1 + \underline{\mathbf{e}}_0 (\mathbf{Q}^*)^{(-1)} \underline{\mathbf{v}}_P} \end{aligned}$$



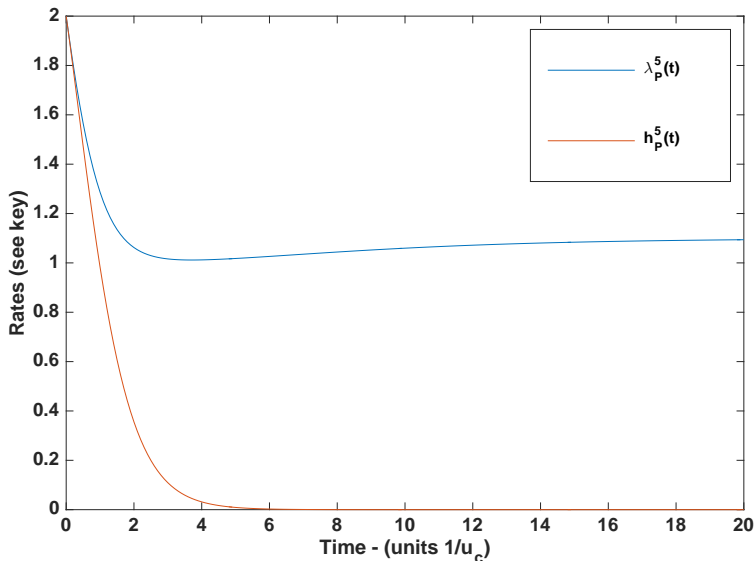


Intuitively, expect pseudogenization rate to go to 0 as  $t \rightarrow \infty$ .  
This is easily proved using fact that

$$\lim_{t \rightarrow \infty} e^{\mathbf{Q}^* t} = 0, \quad (12)$$

Notice

$$\begin{aligned} \lim_{t \rightarrow \infty} h_P(t) &= \frac{\lim_{t \rightarrow \infty} \underline{\mathbf{e}}_0 e^{\mathbf{Q}^* t} \underline{\mathbf{v}}_P}{1 - \lim_{t \rightarrow \infty} \underline{\mathbf{e}}_0 (e^{\mathbf{Q}^* t} - \mathbf{I}) (\mathbf{Q}^*)^{(-1)} \underline{\mathbf{v}}_P} \\ &= \frac{0}{1 + \underline{\mathbf{e}}_0 (\mathbf{Q}^*)^{(-1)} \underline{\mathbf{v}}_P} \\ &= 0. \end{aligned} \quad (13)$$





Some work has been done in the past on modelling this subfunctionalization process.

Some work has been done in the past on modelling this subfunctionalization process.

- Hughes and Liberles (2007,2008) did approximate mechanistic analysis implicitly based on embedded DTMC of process considered here.

Some work has been done in the past on modelling this subfunctionalization process.

- Hughes and Liberles (2007,2008) did approximate mechanistic analysis implicitly based on embedded DTMC of process considered here.
- Later, phenomenological approximations have been used, informed by analysis of Hughes and Liberles (Notable Konrad (2011), Tuefel (2014)).

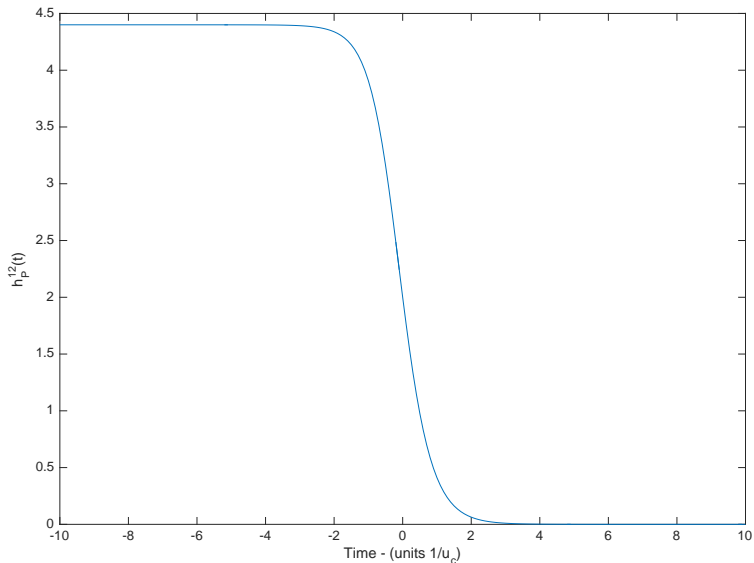
Some work has been done in the past on modelling this subfunctionalization process.

- Hughes and Liberles (2007,2008) did approximate mechanistic analysis implicitly based on embedded DTMC of process considered here.
- Later, phenomenological approximations have been used, informed by analysis of Hughes and Liberles (Notable Konrad (2011), Tuefel (2014)).
- Sigmoid functions have been successful in fitting to real data.

Some work has been done in the past on modelling this subfunctionalization process.

- Hughes and Liberles (2007,2008) did approximate mechanistic analysis implicitly based on embedded DTMC of process considered here.
- Later, phenomenological approximations have been used, informed by analysis of Hughes and Liberles (Notable Konrad (2011), Tuefel (2014)).
- Sigmoid functions have been successful in fitting to real data.

This motivates us to look into the behaviour of the model in negative time!





Limit as  $t \rightarrow -\infty$



UNIVERSITY of  
TASMANIA

Usually interested in physical time. Analysis of negative limit is novel.

Usually interested in physical time. Analysis of negative limit is novel.  
We considered a general CTMC with initial distribution  $\underline{\alpha}$  and structure

$$\mathbf{Q} = \left[ \begin{array}{c|c} \mathbf{Q}^* & \mathbf{V} \\ \hline \mathbf{0} & \mathbf{0} \end{array} \right], \quad (14)$$

Usually interested in physical time. Analysis of negative limit is novel.  
We considered a general CTMC with initial distribution  $\underline{\alpha}$  and structure

$$\mathbf{Q} = \left[ \begin{array}{c|c} \mathbf{Q}^* & \mathbf{V} \\ \hline \mathbf{0} & \mathbf{0} \end{array} \right], \quad (14)$$

$h_P(t; \alpha)$  is an obvious generalization of our earlier function (rate of transition into absorbing state  $P$  assuming not already in  $P$ ).

Usually interested in physical time. Analysis of negative limit is novel.  
We considered a general CTMC with initial distribution  $\underline{\alpha}$  and structure

$$\mathbf{Q} = \left[ \begin{array}{c|c} \mathbf{Q}^* & \mathbf{V} \\ \hline \mathbf{0} & \mathbf{0} \end{array} \right], \quad (14)$$

$h_P(t; \underline{\alpha})$  is an obvious generalization of our earlier function (rate of transition into absorbing state  $P$  assuming not already in  $P$ ).

Using l'Hôpital's rule we get

$$\lim_{t \rightarrow -\infty} h_P(t; \underline{\alpha}) = \lim_{t \rightarrow -\infty} \frac{\underline{\alpha} e^{\mathbf{Q}^* t} \mathbf{Q}^* \mathbf{V}_P}{-\underline{\alpha} e^{\mathbf{Q}^* t} \mathbf{V}_P}. \quad (15)$$

Limit as  $t \rightarrow -\infty$



UNIVERSITY of  
TASMANIA

We diagonalize  $\mathbf{Q}^*$

$$\mathbf{Q}^* = \mathbf{A}^{-1} \Lambda \mathbf{A}$$

We diagonalize  $\mathbf{Q}^*$

$$\mathbf{Q}^* = \mathbf{A}^{-1}\mathbf{\Lambda}\mathbf{A}$$

Then we do some algebra to get

$$\lim_{t \rightarrow -\infty} h_P(t; \underline{\alpha}) = \lim_{t \rightarrow -\infty} \frac{\sum_k [\alpha \mathbf{A}^{-1}]_k e^{\lambda_k t} \lambda_k (\mathbf{A} \mathbf{V}_P)}{-\sum_l [\alpha \mathbf{A}^{-1}]_l e^{\lambda_l t} (\mathbf{A} \mathbf{V}_P)}. \quad (16)$$



We diagonalize  $\mathbf{Q}^*$

$$\mathbf{Q}^* = \mathbf{A}^{-1} \mathbf{\Lambda} \mathbf{A}$$

Then we do some algebra to get

$$\lim_{t \rightarrow -\infty} h_P(t; \underline{\alpha}) = \lim_{t \rightarrow -\infty} \frac{\sum_k [\alpha \mathbf{A}^{-1}]_k e^{\lambda_k t} \lambda_k (\mathbf{A} \mathbf{V}_P)}{-\sum_l [\alpha \mathbf{A}^{-1}]_l e^{\lambda_l t} (\mathbf{A} \mathbf{V}_P)}. \quad (16)$$

Letting  $\lambda_m$  be the eigenvalue of maximum absolute real part of  $\mathbf{Q}^*$  the numerator is



We diagonalize  $\mathbf{Q}^*$

$$\mathbf{Q}^* = \mathbf{A}^{-1} \mathbf{\Lambda} \mathbf{A}$$

Then we do some algebra to get

$$\lim_{t \rightarrow -\infty} h_P(t; \underline{\alpha}) = \lim_{t \rightarrow -\infty} \frac{\sum_k [\alpha \mathbf{A}^{-1}]_k e^{\lambda_k t} \lambda_k(\mathbf{A} \mathbf{V}_P)}{-\sum_l [\alpha \mathbf{A}^{-1}]_l e^{\lambda_l t} (\mathbf{A} \mathbf{V}_P)}. \quad (16)$$

Letting  $\lambda_m$  be the eigenvalue of maximum absolute real part of  $\mathbf{Q}^*$  the numerator is

Numerator

$$\sum_k [\alpha \mathbf{A}^{-1}]_k e^{(\lambda_k - \lambda_m)t} \lambda_k(\mathbf{A} \mathbf{V}_P)$$

Which in the limit is just  $\lambda_m$





Similarly, the denominator is

Denominator

$$- \sum_l [\alpha \mathbf{A}^{-1}]_l e^{(\lambda_l - \lambda_m)t} \lambda_l (\mathbf{A} \mathbf{V}_p)$$



Similarly, the denominator is

Denominator

$$- \sum_l [\alpha \mathbf{A}^{-1}]_l e^{(\lambda_l - \lambda_m)t} \lambda_l (\mathbf{A} \mathbf{V}_p)$$

Which in the limit is just 1.



Similarly, the denominator is

Denominator

$$- \sum_l [\alpha \mathbf{A}^{-1}]_l e^{(\lambda_l - \lambda_m)t} \lambda_l (\mathbf{A} \mathbf{V}_p)$$

Which in the limit is just 1. So, it turns out that

Result

$$\lim_{t \rightarrow -\infty} h_P(t) = -\lambda_m = Sp(\mathbf{Q}) \quad (17)$$



Similarly, the denominator is

Denominator

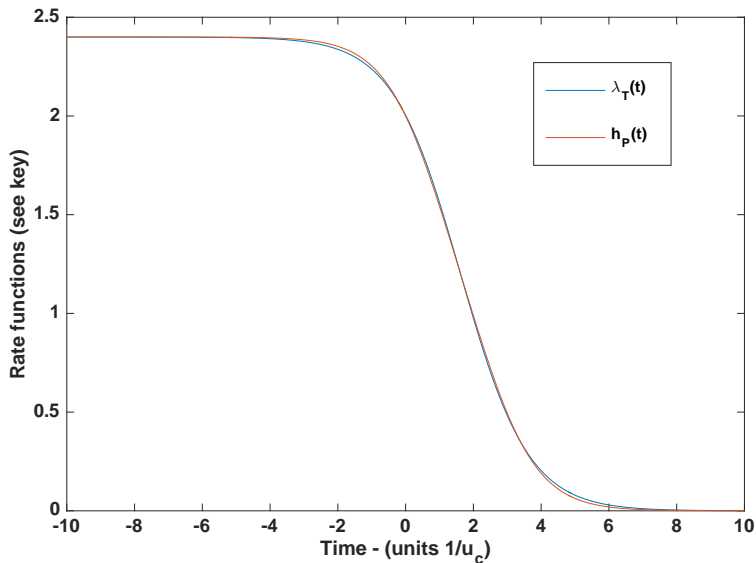
$$- \sum_l [\alpha \mathbf{A}^{-1}]_l e^{(\lambda_l - \lambda_m)t} \lambda_l (\mathbf{A} \mathbf{V}_p)$$

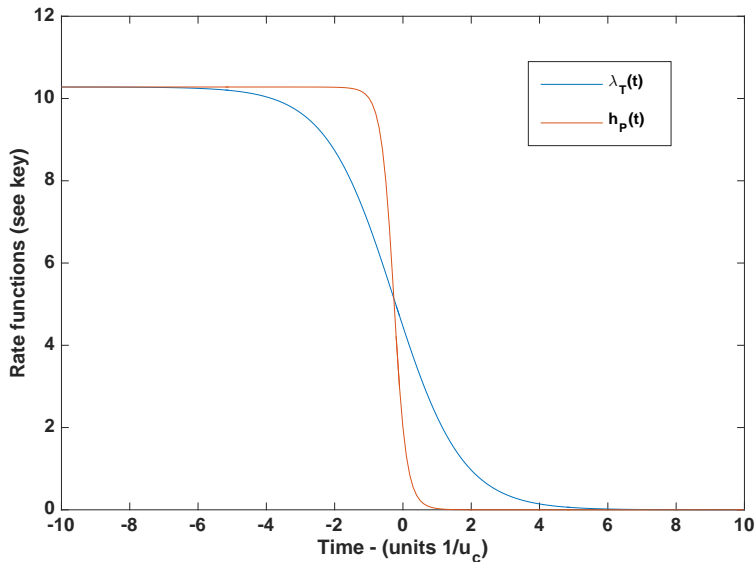
Which in the limit is just 1. So, it turns out that

Result

$$\lim_{t \rightarrow -\infty} h_P(t) = -\lambda_m = Sp(\mathbf{Q}) \quad (17)$$

We can use this to fit the phenomenological approximation of Tuefel et al (2014) to our exact mechanistically function.





Using theory of phase-type distributions we are able to derive and analyze mathematical model for the biological subfunctionalization model assuming only that

Using theory of phase-type distributions we are able to derive and analyze mathematical model for the biological subfunctionalization model assuming only that

- The rate of null mutations in regulatory regions is Poisson  $u_r$



Using theory of phase-type distributions we are able to derive and analyze mathematical model for the biological subfunctionalization model assuming only that

- The rate of null mutations in regulatory regions is Poisson  $u_r$
- The rate of null mutations in the coding region is  $u_c$

Using theory of phase-type distributions we are able to derive and analyze mathematical model for the biological subfunctionalization model assuming only that

- The rate of null mutations in regulatory regions is Poisson  $u_r$
- The rate of null mutations in the coding region is  $u_c$

Pseudogenization rate implied by this model turns out to have the same behaviour as existing phenomenological approximations

Using theory of phase-type distributions we are able to derive and analyze mathematical model for the biological subfunctionalization model assuming only that

- The rate of null mutations in regulatory regions is Poisson  $u_r$
- The rate of null mutations in the coding region is  $u_c$

Pseudogenization rate implied by this model turns out to have the same behaviour as existing phenomenological approximations

- We have provided a means to translate between the exact function and the popular approximation

Using theory of phase-type distributions we are able to derive and analyze mathematical model for the biological subfunctionalization model assuming only that

- The rate of null mutations in regulatory regions is Poisson  $u_r$
- The rate of null mutations in the coding region is  $u_c$

Pseudogenization rate implied by this model turns out to have the same behaviour as existing phenomenological approximations

- We have provided a means to translate between the exact function and the popular approximation
- As well as deriving a host of performance measures

Using theory of phase-type distributions we are able to derive and analyze mathematical model for the biological subfunctionalization model assuming only that

- The rate of null mutations in regulatory regions is Poisson  $u_r$
- The rate of null mutations in the coding region is  $u_c$

Pseudogenization rate implied by this model turns out to have the same behaviour as existing phenomenological approximations

- We have provided a means to translate between the exact function and the popular approximation
- As well as deriving a host of performance measures

This work provides a mathematically rigorous, mechanistically motivated and exact analysis for the fate of gene duplicates.

Using theory of phase-type distributions we are able to derive and analyze mathematical model for the biological subfunctionalization model assuming only that

- The rate of null mutations in regulatory regions is Poisson  $u_r$
- The rate of null mutations in the coding region is  $u_c$

Pseudogenization rate implied by this model turns out to have the same behaviour as existing phenomenological approximations

- We have provided a means to translate between the exact function and the popular approximation
- As well as deriving a host of performance measures

This work provides a mathematically rigorous, mechanistically motivated and exact analysis for the fate of gene duplicates.



Future work will move in two directions

Future work will move in two directions

- Allowing for multiple duplication events to analyze the fates of whole gene families.



Future work will move in two directions

- Allowing for multiple duplication events to analyze the fates of whole gene families.
- Expanding the model to allow for a mixture of sub- and neofunctionalization.