

Annex 2: Algorithms for block event inference, detailed description

Notations:

Input data:

- S is the species tree
- $\{T\}$ is the collection of all the pangenome's gene family trees
- R is a replicon, i.e. a vector of genes of size r

Gene tree indexation:

- G_i is the i -th gene of R gene, and T_i the tree of its gene family
- $t(n_i)$ is the vector of nodes on the lineage of T_i leading to the tip representing G_i , i.e. the path between the gene tree root and the tip
- $n_{i,j}$ is the j -th node on t_i , counting from the tip representing G_i (that is $n_{i,0}$), i.e. $n_{i,j} = t(n_i)[j]$

ODT event indexation:

- gene tree nodes are annotated with events, and given the event collection $\{E\}$, the map Ev links a node n with the event $E_n = Ev(n, \{E\})$
- $E_{i,j}$ is the event associated to $n_{i,j}$
- m is a node in S representing an ancestral species
- $L(X) = \{m\}$ is the location of the event or block event X , characterized by the set of possible nodes in S where it can have occurred

Block classes:

- B is a leaf block event, instance of class *LBlock*, characterized by the following attributes: the ordered vector of genes it gathers $(G)_B$, the collection of associated events $\{E\}_B$, and its aggregate location $L(B)$
- A is an ancestral block event, instance of class *ABlock*, characterized by the following attributes: the collection of leaf block events it gathers $(B)_A$, the collection of associated gene tree events $\{E\}_A$, and its aggregate location $L(A)$

Block construction (parameters):

- g is the maximum allowed gap length for leaf block event construction

Block construction (dynamic variables):

- $(A)_c$ and $(B)_c$ are the current lists of all ancestral and leaf block events, respectively
- $\{A\}_E$ is the set of ancestral blocks that are currently linked to a gene tree event E (through any of their member leaf blocks)
- $getanc$ is a function such as $\{A\}_E = getanc(E)$
- P_B is the list of ancestral blocks with which leaf block B shares gene tree events, and could be putatively assigned to
- $A(B)$ is the ancestral block event to which leaf block B is currently assigned
- c is the currently opened gap size (during leaf block event construction)

Helper Functions:

```
// instantiate new blocks
def new(LBlock, E, G) :
    set B ;
    set L(B) = L(E) ;
    set (G)B = (G) ;
    set {E}B = {E} ;
    return B ;

def new(ABlock, B) :
    set A ;
    set (B)A = (B) ;
    set L(A) = L(B) ;
    set {E}A = {E}B ;
    return A ;

// add event or leaf block to leaf or ancestral block, respectively
def add(B, E, G):
    (G)B += (G) ;
    {E}B = {E}B ∪ {E} ;
    L(B) = L(B) ∩ L(E) ;

def add(A, B):
    (B)A += (B) ;
    {E}A = {E}A ∪ {E}B ;
    L(A) = L(A) ∩ L(B) ;

// test compatibility between two instances of the following:
// events, leaf blocks or ancestral blocks
def compat(X, Y):
    Lcompa = L(X) ∩ L(Y) ;
    return (Lcompa ≠ ∅) ;
```

```
// split block B before element i
def split(B, i):
    // create a new block B' covering the right-hand part of B
    B' = new(LBlock, {E}B[i:], (G)B[i:])
    // restrict block B to its left-hand part
    set (G)B = (G)B[1:i-1] ;
    set {E}B = {E}B[1:i-1] ;
    set L(B) = L(E1) ∩ L(E2) ∩ ... ∩ L(Ei-1)
    return ( B, B' )

def merge(A1, A2):
    (B)A1 = (B)A1 ∪ (B)A2 ;
    L(A1) = L(A1) ∩ L(A2) ;
    {E}A += {E}B ;
    del A2 ;
    return A1 ;

// obtain all ancestral blocks linked to
// (i.e. currently known to include) a gene tree event E
def getanc(E, (A)c) :
    set {A}E = {}
    for A in (A)c ; do
        if (( E ∈ {E}A )) ; then {A}E += (A) ; endif
    end ;
    return {A}E ;
```

```
// resolve coordinates incompatibilities between leaf blocks
// and linked ancestral blocks by splitting the leaf block into
// as many needed leaf blocks that are compatible with their
// respective linked ancestral blocks
def resolve(B, A):
    I = () ; // vector of indexes of events in B not compatible with A
    for i in (i >= 1; i <= len({E}B; i++) ; do
        Ei = {E}B[i]
        if ( not compat(Ei, A) ) ; then
            I += (i) ;
        endif
    end
    for i in I ; do
        if ( i > 1 ) ; then
            // split block B before incompatible event Ei
            (B, B') = split(B, i) ;
            // recursive call ; returns a vector of blocks
            // (compatible part of the original blocks always first)
            return ( B ) + resolve(B', A)
        else
            // split block at the next compatible event
            for (k >= 2; k <= len({E}B; k++) ; do
                if ( k ∉ I ) ; then
                    (B, B') = split(B, k) ;
                    break // the for k loop
                endif
            end
            // recursive call ; returns a vector of blocks
            return resolve(B', A) + ( B )
        endif
    end
    // end of recursive call
    return ( B )
```

Algorithm 1: Construction of leaf blocks.**INPUT:**

- a vector of r genes G_i ordered on a replicon R (single DNA molecule in an extant genome),
- the collection $\{T\}$ of all reconciled gene trees, of which every node n from the set $\{n\}_T$ is annotated with an event E characterized by a location $L(E)$ in the species tree,
- a map F associating a gene G_i to its respective reconciled gene family tree $T_i = F(G_i, \{T\})$

OUTPUT:

- a list of leaf block events $(B)_R$ for replicon R , and their respective attributes

// start with an empty list of leaf block events

$(B)_R = ()$

for $((i \geq 1 ; i \leq r ; i++))$; **do**

 // access the i -th gene on R

$G_i = R[i]$;

 // access the corresponding gene tree, the tip node representing the gene and the lineage above

$T_i = F(G_i, \{T\})$; $n_i = T_i[G_i]$; $t_i = t(n_i)$;

 // explore the lineage above

for $((j \geq 1 ; n_{i,j} \neq \text{root}(T_i) ; j++))$; **do**

$n_{i,j} = t_i[j]$

$E_{i,j} = \text{Ev}(n_{i,j}, \{E\})$

 // instantiate a new block event B with $E_{i,j}$ as seed

$B = \text{new}(\text{LBlock}, E_{i,j}, G_i)$;

$c = 0$;

$k = i + 1$;

while $((c \leq g \text{ and } k < r))$; **do**

 // scan the right-hand neighbour genes G_k on the replicon

$G_k = R[k]$;

$T_k = F(G_k, \{T\})$; $n_k = T_k[G_k]$; $t_k = t(n_k)$;

$l = 1$;

while $((l \leq \text{length}(t_k)))$; **do**

 // lineage above neighbour gene on the right G_k is explored

$n_{k,l} = t_k[l]$; $E_{k,l} = \text{Ev}(n_{k,l}, \{E\})$;

 // test compatibility of events locations, i.e. their non-null intersection

if $(\text{compat}(B, E_{k,l}))$; **then**

$\text{add}(B, E_{k,l}, G_k)_B$;

$c = 0$;

break // stop the exploration of T_k for events compatible to $E_{i,j}$

endif

if $((l == \text{length}(t_k)))$; **then**

 // no compatible event was found in lineage above G_k

$(G)_B += (G_k)$; // add G_k to B as a gap gene (no refinement of $L(B)$)

 // increment gap count

$c++$;

endif

$l++$;

end

$k++$

end

if $((c == g))$; **then**

 // there is no more compatible neighbour to add, the trail of gap genes must be trimmed from the block

$(G)_B = (G)_B[1:\text{len}(B)-g]$

$(B)_R += (B)$

endif

end

return $(B)_R$

// end of the block event reconstruction for replicon R ; repeat for all replicons in the genome dataset

Algorithm 2: Construction of ancestral blocks.**INPUT:**

- the list of all leaf blocks $(B)_c$

OUTPUT:

- the list of ancestral blocks $(A)_c$
- added attributes to leaf blocks, $A(B)$ pointing to their respective ancestral blocks

// initiate the ancestral block list

$(A)_c = ()$;

for $(b \geq 1 ; b \leq \text{len}((B)_c) ; b++)$ **do** ;

$B = (B)_c[b]$;

// search for already existing ancestral blocks linked to gene tree events from B

$P_B = ()$;

set $A(B) = \emptyset$;

for E in $\{E\}_B$; **do**

$\{A\}E = \text{getanc}(E, (A)_c) ; P_B += \{A\}E$;

end

if $(P_B \neq \emptyset)$; **then**

$A = \text{new}(\text{ABlock}, B) ; P_B += (A)$;

$(A)_c += (A)$;

endif

for $(i \geq 1 ; i \leq \text{len}(P_B) ; i++)$; **do**

$A_i = P_B[i]$

// test compatibility of events locations, i.e. their non-null intersection

if $(\text{not } \text{compat}(B, A_i))$; **then**

// split block into compatible and incompatible parts

// compatible part of the original blocks is always

// returned as first element of the vector

$(B, B', \dots) = \text{resolve}(B, A_i)$;

// append new block to global list for further allocation

$(B)_c += (B', \dots)$;

endif

if $(A(B) == \emptyset)$; **then**

$A(B) = A_i$;

else

$A(B) = \text{merge}(A_i, A_B)$;

endif

end

end