

Reconstructing Cancer Karyotypes from paired end reads and CN data using ILP

Rami Eitan and Ron Shamir, Tel-Aviv University

16/12/15

• The current understanding of cancer suggests that it is a disease driven by somatic mutations that accumulate in the genome, within a certain tissue, during the lifetime of an individual.

- The current understanding of cancer suggests that it is a disease driven by somatic mutations that accumulate in the genome, within a certain tissue, during the lifetime of an individual.
- We consider 3 types of intra-chromosomal rearrangements:

- The current understanding of cancer suggests that it is a disease driven by somatic mutations that accumulate in the genome, within a certain tissue, during the lifetime of an individual.
- We consider 3 types of intra-chromosomal rearrangements:



- The current understanding of cancer suggests that it is a disease driven by somatic mutations that accumulate in the genome, within a certain tissue, during the lifetime of an individual.
- We consider 3 types of intra-chromosomal rearrangements:



- The current understanding of cancer suggests that it is a disease driven by somatic mutations that accumulate in the genome, within a certain tissue, during the lifetime of an individual.
- We consider 3 types of intra-chromosomal rearrangements:



- The current understanding of cancer suggests that it is a disease driven by somatic mutations that accumulate in the genome, within a certain tissue, during the lifetime of an individual.
- We consider 3 types of intra-chromosomal rearrangements:



• In addition we also consider the inter-choromosomal rearrangement of translocation.

• In addition we also consider the inter-choromosomal rearrangement of translocation.



Balanced reciprocal

• In addition we also consider the inter-choromosomal rearrangement of translocation.



• In addition we also consider the inter-choromosomal rearrangement of translocation.



• Over time the rearrangements accumulate and the cancer cell karyoptype differs more from its healthy counterpart



• CGH array – Detect Copy Number Variations (CNV), resolution as low as 100 Kilo bases.

 CGH array – Detect Copy Number Variations (CNV), resolution as low as 100 Kilo bases.



• Paired Ends Reads – Detect novel adjacencies in the genome compared to the reference.

- Paired Ends Reads Detect novel adjacencies in the genome compared to the reference.
 - Sample DNA sequence *S* is cut into small fragments (200-500 bp)

- Paired Ends Reads Detect novel adjacencies in the genome compared to the reference.
 - Sample DNA sequence *S* is cut into small fragments (200-500 bp)
 - Each end of the fragment (36 bp) is then aligned against a reference genome *R*.

- Paired Ends Reads Detect novel adjacencies in the genome compared to the reference.
 - Sample DNA sequence *S* is cut into small fragments (200-500 bp)
 - Each end of the fragment (36 bp) is then aligned against a reference genome *R*.
 - Concordant reads both ends aligned to the known expected distance and orientation.



Concordant read

- Paired Ends Reads Detect novel adjacencies in the genome compared to the reference.
 - Sample DNA sequence *S* is cut into small fragments (200-500 bp)
 - Each end of the fragment (36 bp) is then aligned against a reference genome *R*.
 - Concordant reads both ends aligned to the known expected distance and orientation.
 - Discordant reads are mapped to different locations on the reference genome or with an unexpected relative orientation





Concordant read

Discordant read





Tandem duplication





Deletion









Deletion



Tandem duplication



Inversion



• We have a *reference* (normal) genome and an unknown *target* genome.

- We have a *reference* (normal) genome and an unknown *target* genome.
- Breakpoints divide each chromosome into segments.

- We have a *reference* (normal) genome and an unknown *target* genome.
- Breakpoints divide each chromosome into segments.
- We call the end points of a segment I its *tail* and *head*. $I = [t_I, h_I]$

- We have a *reference* (normal) genome and an unknown *target* genome.
- Breakpoints divide each chromosome into segments.
- We call the end points of a segment I its *tail* and *head*. $I = [t_I, h_I]$
- A *bridge* is a connection between two segment end points that are adjacent in the target genome but not in the reference.

- We have a *reference* (normal) genome and an unknown *target* genome.
- Breakpoints divide each chromosome into segments.
- We call the end points of a segment I its *tail* and *head*. $I = [t_I, h_I]$
- A *bridge* is a connection between two segment end points that are adjacent in the target genome but not in the reference.



The input is:

The input is:

• A sequence of intervals representing the reference genome $\mathcal{I} = \{I_1, \dots, I_n\}.$

The input is:

- A sequence of intervals representing the reference genome
 \$\mathcal{I} = \{I_1, \ldots, I_n\}\$.
- The copy number profile of the intervals, each Interval I_j has CN N_j .



The input is:

- A sequence of intervals representing the reference genome
 \$\mathcal{I} = \{I_1, \ldots, I_n\}\$.
- The copy number profile of the intervals, each Interval I_j has CN N_j .
- The set of bridges $\{a_i, b_i\}_{i=1}^m$ and the support μ_i for each bridge.




• weighted graph G(V, E, w) whose vertices are the interval extremities.







- weighted graph G(V, E, w) whose vertices are the interval extremities.
- Three types of edges:
 - *interval edge* $e_I(t_i, h_i) \in E_I$ for each interval $I_i = [t_i, h_i]$, weighted N_i







- weighted graph G(V, E, w) whose vertices are the interval extremities.
- Three types of edges:
 - *interval edge* $e_I(t_i, h_i) \in E_I$ for each interval $I_i = [t_i, h_i]$, weighted N_i
 - reference edge $e_R(h_i, t_{i+1}) \in E_R$, connecting the head if an interval to the tail of the one adjacent to it in the reference.







- weighted graph G(V, E, w) whose vertices are the interval extremities.
- Three types of edges:
 - *interval edge* $e_I(t_i, h_i) \in E_I$ for each interval $I_i = [t_i, h_i]$, weighted N_i
 - reference edge $e_R(h_i, t_{i+1}) \in E_R$, connecting the head if an interval to the tail of the one adjacent to it in the reference.
 - bridge (or variant) edge $e_V(a_i, b_j) \in E_V$ connecting to extremities that are adjacent in the target.







Bridge graph cont.

• Transform all undirected edges to antiparallel directed edges.





Bridge graph cont.

- Transform all undirected edges to antiparallel directed edges.
- Add weights to bridges, representing support score.





Bridge graph cont.

- Transform all undirected edges to antiparallel directed edges.
- Add weights to bridges, representing support score.
- Denote S ⊆ V the *telomere nodes* that start or end each reference chromosome.





• Given a bridge graph G(V, E, w), a valid path p is a path through G that:

- Given a bridge graph G(V, E, w), a valid path p is a path through G that:
 - Start and ends on a telomeric node

- Given a bridge graph G(V, E, w), a valid path p is a path through G that:
 - Start and ends on a telomeric node
 - Alternates between interval and non-interval edges.

- Given a bridge graph G(V, E, w), a valid path p is a path through G that:
 - Start and ends on a telomeric node
 - Alternates between interval and non-interval edges.



- Given a bridge graph G(V, E, w), a valid path p is a path through G that:
 - Start and ends on a telomeric node
 - Alternates between interval and non-interval edges.
- Formally:

- Given a bridge graph G(V, E, w), a valid path p is a path through G that:
 - Start and ends on a telomeric node
 - Alternates between interval and non-interval edges.
- Formally:
 - denote $f_p(e_i)$ the number of times edge e_i is traversed in either direction.

- Given a bridge graph G(V, E, w), a valid path p is a path through G that:
 - Start and ends on a telomeric node
 - Alternates between interval and non-interval edges.
- Formally:
 - denote $f_p(e_i)$ the number of times edge e_i is traversed in either direction.
 - $E_{I_{\leftarrow}}(v), E_{I \rightarrow}(v), E_{R \leftarrow}(v), E_{R \rightarrow}(v), E_{V \leftarrow}(v), E_{V \rightarrow}(v)$ the set of interval, reference and variant edges that go in and out of v respectively.

- Given a bridge graph G(V, E, w), a valid path p is a path through G that:
 - Start and ends on a telomeric node
 - Alternates between interval and non-interval edges.
- Formally:
 - denote $f_p(e_i)$ the number of times edge e_i is traversed in either direction.
 - $E_{I_{\leftarrow}}(v), E_{I \rightarrow}(v), E_{R \leftarrow}(v), E_{R \rightarrow}(v), E_{V \leftarrow}(v), E_{V \rightarrow}(v)$ the set of interval, reference and variant edges that go in and out of v respectively.

•
$$f_P\left(E_{I_{\rightarrow}}(v)\right) = f_P\left(E_{R_{\leftarrow}}(v)\right) + f_P\left(E_{V_{\leftarrow}}(v)\right)$$

 $\forall_{v \notin S}$

- Given a bridge graph G(V, E, w), a valid path p is a path through G that:
 - Start and ends on a telomeric node
 - Alternates between interval and non-interval edges.
- Formally:
 - denote $f_p(e_i)$ the number of times edge e_i is traversed in either direction.
 - $E_{I_{\leftarrow}}(v), E_{I \rightarrow}(v), E_{R \leftarrow}(v), E_{R \rightarrow}(v), E_{V \leftarrow}(v), E_{V \rightarrow}(v)$ the set of interval, reference and variant edges that go in and out of v respectively.

•
$$f_P\left(E_{I_{\rightarrow}}(v)\right) = f_P\left(E_{R_{\leftarrow}}(v)\right) + f_P\left(E_{V_{\leftarrow}}(v)\right)$$

 $\forall_{v\notin S}$
• $f_P\left(E_{I_{\leftarrow}}(v)\right) = f_P\left(E_{R_{\rightarrow}}(v)\right) + f_P\left(E_{V_{\rightarrow}}(v)\right)$
 $\forall_{v\notin S}$

- Given a bridge graph G(V, E, w), a valid path p is a path through G that:
 - Start and ends on a telomeric node
 - Alternates between interval and non-interval edges.
- Formally:
 - denote $f_p(e_i)$ the number of times edge e_i is traversed in either direction.
 - $E_{I,-}(v), E_{I\to}(v), E_{R\leftarrow}(v), E_{R\to}(v), E_{V\leftarrow}(v), E_{V\to}(v)$ the set of interval, reference and variant edges that go in and out of v respectively.

•
$$f_P\left(E_{I_{\rightarrow}}(v)\right) = f_P\left(E_{R_{\leftarrow}}(v)\right) + f_P\left(E_{V_{\leftarrow}}(v)\right)$$

 $\forall_{v \notin S}$
• $f_P\left(E_{I_{\leftarrow}}(v)\right) = f_P\left(E_{R_{\rightarrow}}(v)\right) + f_P\left(E_{V_{\rightarrow}}(v)\right)$
 $\forall_{v \notin S}$

• $f_P(e) \in \mathbb{N}^0$

Distance of path from observed data

• We define a *discordancy score* to measure how consistent a path is with the observed data.

Distance of path from observed data

• We define a *discordancy score* to measure how consistent a path is with the observed data.

$$d_G(f_P) = \sum_{e \in E_I} \frac{l_e}{L} |f_P(e) - w(e)| + \alpha \sum_{\substack{e \in E_V \\ e \notin P}} \frac{w(e)}{\mu}$$

Distance of path from observed data

• We define a *discordancy score* to measure how consistent a path is with the observed data.

$$d_G(f_P) = \sum_{e \in E_I} \frac{l_e}{L} |f_P(e) - w(e)| + \alpha \sum_{\substack{e \in E_V \\ e \notin P}} \frac{w(e)}{\mu}$$

- w(e) the weight of an edge e
- l_e length of the interval represented by e
- L Total length of all intervals
- $\mu = \sum_{e \in E_V} w(e)$ the sum of the support score for all bridges

ILP formulation

- We want to find a path that has the lowest discordancy score.
- This can be formulated as an ILP formulation:

ILP formulation

- We want to find a path that has the lowest discordancy score.
- This can be formulated as an ILP formulation:
- Minimize:

$$d_G(f_P) = \sum_{e \in E_I} \frac{l_e}{L} \left| x_{e \to}^I + x_{e \leftarrow}^I - w(e) \right| + \alpha \sum_{e \in E_V} \frac{w(e)}{\mu} \left(1 - \min\left(1, x_{e \to}^B + x_{\leftarrow}^B\right) \right)$$

ILP formulation

- We want to find a path that has the lowest discordancy score.
- This can be formulated as an ILP formulation:
- Minimize:

$$d_G(f_P) = \sum_{e \in E_I} \frac{l_e}{L} \left| x_{e \to}^I + x_{e \leftarrow}^I - w(e) \right| + \alpha \sum_{e \in E_V} \frac{w(e)}{\mu} \left(1 - \min\left(1, x_{e \to}^B + x_{\leftarrow}^B\right) \right)$$

- Subject to:
 - $\forall_i x_i \in \mathbb{N}^0$

•
$$\forall_{v \notin S} \sum_{e_i \in E_{I \to}(v)} x_{i \to}^I = \sum_{e_i \in E_{R \leftarrow}(v)} x_{i \leftarrow}^R + \sum_{e_i \in E_{V \leftarrow}(v)} x_{i \leftarrow}^R$$

•
$$\forall_{v \notin S} \sum_{e_i \in E_{I_{\leftarrow}}(v)} x_{i_{\leftarrow}}^I = \sum_{e_i \in E_{R_{\rightarrow}}(v)} x_{i_{\rightarrow}}^R + \sum_{e_i \in E_{V \rightarrow}(v)} x_{i_{\rightarrow}}^R$$

• Simulate a rearranged karyotype .

- Simulate a rearranged karyotype .
- Apply noise

- Simulate a rearranged karyotype .
- Apply noise
 - For CN: Add $x \sim N(0, \epsilon)$

- Simulate a rearranged karyotype .
- Apply noise
 - For CN: Add $x \sim N(0, \epsilon)$
- For support score: Draw from $Exp(\lambda)$

- Simulate a rearranged karyotype .
- Apply noise
 - For CN: Add $x \sim N(0, \epsilon)$
- For support score: Draw from $Exp(\lambda)$
- "Miss" bridges with a probability p for each bridge.

• C - The number of chromosomes (default: 5).

- *C* The number of chromosomes (default: 5).
- *N* The number of structural and numerical operations applied (default: 5).

- *C* The number of chromosomes (default: 5).
- *N* The number of structural and numerical operations applied (default: 5).
- ϵ The standard deviation of the noise in the CN profile data (default: 0.2)

- *C* The number of chromosomes (default: 5).
- *N* The number of structural and numerical operations applied (default: 5).
- ϵ The standard deviation of the noise in the CN profile data (default: 0.2)
- *p* The probability to completely miss a bridge (default: 0.05).

Simulations – Correctness measures
Let T be the simulated (true) karyotype, let T^* be the simulated noisy karyotype, and let S be the karyotype produced by the algorithm.

• Correct - S and T have the same CN profile and use the same bridges.

- Correct S and T have the same CN profile and use the same bridges.
- Equal Copy Number (ECN) *S* and *T* have the same CN profile.

- Correct S and T have the same CN profile and use the same bridges.
- Equal Copy Number (ECN) S and T have the same CN profile.
- Equal or Better Score (EBS) S is closer to T^{*} than to T.

- Correct S and T have the same CN profile and use the same bridges.
- Equal Copy Number (ECN) S and T have the same CN profile.
- Equal or Better Score (EBS) S is closer to T^* than to T.
- Equivalent for Observed Bridges (EOB) S is equivalent to T if we ignore the missing bridges.

- Correct S and T have the same CN profile and use the same bridges.
- Equal Copy Number (ECN) S and T have the same CN profile.
- Equal or Better Score (EBS) S is closer to T^* than to T.
- Equivalent for Observed Bridges (EOB) S is equivalent to T if we ignore the missing bridges.
- CN score The fraction of the intervals that have the correct copy number.



Mean success rate for the base scenario

The effect of bridge support weight in the objective

Preformance as a function of alpha





The effect of the number of operations



Number of operations



The effect of the number of chromosomes

Results as a function of the probability to miss a bridge



Probability to miss a bridge



Real data: Malhorta et al. (2013)





Real data: GBM 10





Real data: GBM 10



Real data: LUAD 6, chromosome 1



Real data: LUAD 6, chromosome 1



Real data: LUSC 5



Real data: LUSC 5



Real data: LUSC 5











