

Data and text mining

# CAPRI: efficient inference of cancer progression models from cross-sectional data

Daniele Ramazzotti<sup>1,\*</sup>, Giulio Caravagna<sup>1</sup>, Loes Olde Loohuis<sup>2</sup>,  
Alex Graudenzi<sup>1</sup>, Ilya Korsunsky<sup>3</sup>, Giancarlo Mauri<sup>1,4</sup>, Marco Antoniotti<sup>1</sup>  
and Bud Mishra<sup>3</sup>

<sup>1</sup>Department of Informatics, Systems and Communication, University of Milan-Bicocca, Milan, Italy, <sup>2</sup>Center for Neurobehavioral Genetics, University of California Los Angeles, Los Angeles, CA, USA, <sup>3</sup>Courant Institute of Mathematical Sciences, New York University, New York, NY, USA and <sup>4</sup>SYSBIO Centre of Systems Biology, Milano, Italy

\*To whom correspondence should be addressed.  
Associate Editor: Jonathan Wren

Received on January 19, 2015; revised on April 7, 2015; accepted on May 4, 2015

## Abstract

**Summary:** We devise a novel inference algorithm to effectively solve the *cancer progression model reconstruction* problem. Our empirical analysis of the accuracy and convergence rate of our algorithm, *CAnceR PRogression Inference* (CAPRI), shows that it outperforms the state-of-the-art algorithms addressing similar problems.

**Motivation:** Several cancer-related genomic data have become available (e.g. *The Cancer Genome Atlas*, TCGA) typically involving hundreds of patients. At present, most of these data are aggregated in a *cross-sectional* fashion providing all measurements at the time of diagnosis. Our goal is to infer cancer ‘progression’ models from such data. These models are represented as directed acyclic graphs (DAGs) of collections of ‘selectivity’ relations, where a mutation in a gene *A* ‘selects’ for a later mutation in a gene *B*. Gaining insight into the structure of such progressions has the potential to improve both the stratification of patients and personalized therapy choices.

**Results:** The CAPRI algorithm relies on a scoring method based on a *probabilistic theory* developed by Suppes, coupled with *bootstrap* and *maximum likelihood* inference. The resulting algorithm is efficient, achieves high accuracy and has good complexity, also, in terms of convergence properties. CAPRI performs especially well in the presence of noise in the data, and with limited sample sizes. Moreover CAPRI, in contrast to other approaches, robustly reconstructs different types of confluent trajectories despite irregularities in the data. We also report on an ongoing investigation using CAPRI to study *atypical Chronic Myeloid Leukemia*, in which we uncovered non trivial selectivity relations and exclusivity patterns among key genomic events.

**Availability and implementation:** CAPRI is part of the *TRanslational ONCOlogy* R package and is freely available on the web at: <http://bimib.disco.unimib.it/index.php/Tronco>

**Contact:** [daniele.ramazzotti@disco.unimib.it](mailto:daniele.ramazzotti@disco.unimib.it)

**Supplementary information:** [Supplementary data](#) are available at *Bioinformatics* online.

## 1 Introduction

Analysis and interpretation of the fast-growing biological data sets that are currently being curated from laboratories all over the world require sophisticated computational and statistical methods.

Motivated by the availability of genetic patient data, we focus on the problem of *reconstructing progression models* of cancer. In particular, we aim to infer the plausible sequences of *genomic alterations* that, by a process of *accumulation*, selectively make a tumor fitter to survive, expand and diffuse (i.e. metastasize). Along the trajectories of progression, a tumor (monotonically) acquires or ‘activates’ mutations in the genome, which, in turn, produce progressively more ‘viable’ clonal subpopulations over the so-called *cancer evolutionary landscape* (cf., Merlo *et al.*, 2006; Huang *et al.*, 2009; Vogelstein *et al.*, 2013).

Knowledge of such progression models is very important for drug development and in therapeutic decisions. For example, it has been known that for the same cancer type, patients in different stages of different progressions respond differently to different treatments.

Several datasets are currently available that aggregate diverse cancer-patient data and report in-depth mutational profiles, including e.g. structural changes (e.g. inversions, translocations, copy-number variations) or somatic mutations (e.g. point mutations, insertions, deletions, etc.). An example of such a dataset is *The Cancer Genome Atlas* (TCGA) (cf., NCI and the NHGRI, 2005). These data, by their very nature, only give a snapshot of a given tumor sample, mostly from biopsies of untreated tumor samples at the time of diagnoses. It still remains impractical to track the tumor progression in any single patient over time, thus limiting most analysis methods to work with *cross-sectional* data. (Unlike longitudinal studies, these cross-sectional data are derived from samples that are collected at unknown time points, and can be considered as ‘static’.)

To rephrase, we focus on the problem of *cancer progression models reconstruction from cross-sectional data*. The problem is not new and, to the best of our knowledge, two threads of research starting in the late 90s have addressed it. The first category of works examined mostly gene-expression data to reconstruct the temporal ordering of samples (cf., Magwene *et al.*, 2003; Gupta and Bar-Joseph, 2008). The second category of works looked at inferring cancer progression models of increasing model-complexity, starting from the simplest tree models (cf. Desper *et al.*, 1999) to more complex graph models (cf., Gerstung *et al.*, 2009); see the next subsection for an overview of the state of the art. Building on our previous work described in Olde Loohuis *et al.* (2014) we present a novel and comprehensive algorithm of the second category that addresses this problem.

The new algorithm proposed here is called *CAnceR PRogression Inference* (CAPRI) and is part of the *TRanslational ONCOlogy* (TRONCO) package (cf., Antoniotti *et al.*, 2014). Starting from cross-sectional genomic data, CAPRI reconstructs a probabilistic progression model by inferring ‘selectivity relations’, where a mutation in a gene *A* ‘selects’ for a later mutation in a gene *B*. These relations are depicted in a combinatorial graph and resemble the way a mutation exploits its ‘selective advantage’ to allow its host cells to expand clonally. Among other things, a selectivity relation implies a putatively invariant temporal structure among the genomic alterations (i.e. *events*) in a specific cancer type. In addition, these relations are expected to also imply ‘probability raising’ for a pair of events in the following sense: Namely, a selectivity relation between a pair of events here signifies that the presence of the earlier genomic alteration (i.e. the *upstream event*) that is advantageous in a

Darwinian competition scenario increases the probability with which a subsequent advantageous genomic alteration (i.e. the *downstream event*) appears in the clonal evolution of the tumor. Thus the selectivity relation captures the effects of the evolutionary processes, and not just correlations among the events and imputed clocks associated with them. As an example, we show in (Fig. 1) the selectivity relation connecting a mutation of *EGFR* to the mutation of *CDK*.

Consequently, an inferred selectivity relation suggests mutational profiles in which certain samples (early-stage patients) display specific alterations only (e.g. the alteration characterizing the beginning of the progression), while certain other samples (e.g. late-stage patients) display a superset subsuming the early mutations (as well as alterations that occur subsequently in the progression).

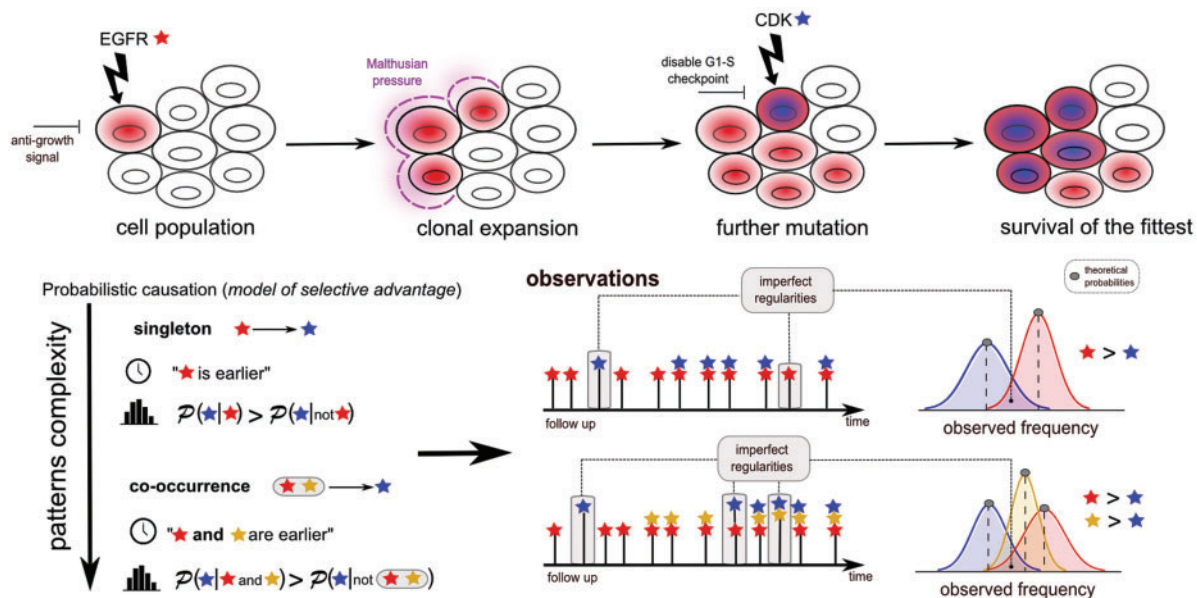
Various kinds of genomic aberrations are suitable as input data, and include somatic point/indel mutations, copy-number alterations, etc., provided that they are *persistent*, i.e. once an alteration is acquired no other genomic event can restore the cell to the non-mutated (i.e. *wild type*) condition. (For instance, epigenetic alterations such as methylation and alterations in gene expression are not directly usable as input data for the algorithm. Notice that the selection of the relevant events is beyond the scope of this work and requires a further upstream pipeline, such as that provided, for instance, in Tamborero *et al.*, 2013; Vogelstein *et al.*, 2013.)

The selectivity relations that CAPRI reconstructs are ranked and subsequently further refined by means of a hybrid algorithm, which reasons over time, mechanism and chance, as follows. CAPRI’s overall scoring methods combine topological constraints grounded on Patrick Suppes’ conditions of probabilistic causation (see e.g. Suppes, 1970), with a *maximum likelihood-fit* procedure (cf., Koller and Friedman, 2009) and derives much of its statistical power from the application of *bootstrap* procedures (see e.g. Efron, 1982). CAPRI returns a *graphical model* of a complex selectivity relation among events which captures the essential aspects of cancer evolution: branches, confluences and independent progressions. In the specific case of confluences, CAPRI’s ability to infer them is related to the complexity of the ‘patterns’ they exhibit, expressed in a logical fashion. As pointed out by other approaches (cf., Beerenwinkel *et al.*, 2007), this strategy requires trading off complexity for expressivity of the inferred models, and results in two execution modes for the algorithm: supervised and unsupervised, which we discuss in detail in Sections 2 and 3.

In Section 3 we show that CAPRI enjoys a set of attractive properties in terms of its complexity, soundness and expressivity, even in the presence of uniform *noise* in the input data—e.g. due to *genetic heterogeneity* and experimental errors. Although many other approaches enjoy similar asymptotic properties, we show that CAPRI can compute accurate results with surprisingly small sample sizes (cf., Section 4). Moreover, to the best of our knowledge, based on extensive synthetic data simulations, CAPRI outperforms all the competing procedures with respect to all desirable performance metrics. We conclude by showing an application of CAPRI to reconstruct a progression model for *atypical Chronic Myeloid Leukemia* (aCML) using a recent exome sequencing dataset, first presented in Piazza *et al.* (2013).

### 1.1 State of the art

For an extensive review on *cancer progression model reconstruction* we refer to the recent survey by Beerenwinkel *et al.* (2014). In brief, progression models for cancer have been studied starting with the seminal work of Vogelstein *et al.* (1988) where, for the first time, cancer progression was described in terms of a directed path by assuming the existence of a unique and most likely temporal order



**Fig. 1.** Selectivity relation in tumor evolution. The *CAnceR PRogression Inference* (CAPRI) algorithm examines cancer patients' genomic cross-sectional data to determine relationships among genomic alterations (e.g. somatic mutations, copy-number variations, etc.) that modulate the somatic evolution of a tumor. When CAPRI concludes that aberration *a* (say, an EGFR mutation) 'selects for' aberration *b* (say, a CDK mutation), such relations can be rigorously expressed using Suppes' conditions, which postulates that if *a* selects *b*, then *a* occurs before *b* (*temporal priority*) and occurrences of *a* raises the probability of emergence of *b* (*probability raising*). Moreover, CAPRI is capable of reconstructing relations among more complex boolean combination of events, as shown in the bottom panel and discussed in the Approach section

of genetic mutations. Vogelstein *et al.* (1988) manually created a (colorectal) cancer progression from a genetic and clinical point of view. More rigorous and complex algorithmic and statistical automated approaches have appeared subsequently. As stated already, the earliest thread of research simply sought more generic progression models that could assume tree-like structures. The *oncogenetic tree model* captured evolutionary branches of mutations (cf., Desper *et al.*, 1999; Szabo and Boucher, 2002) by optimizing a correlation-based score. Another popular approach to reconstruct tree structures appears in Desper *et al.* (2000). Other general Markov chain models such as, e.g. Hjelm *et al.* (2006) reconstruct more flexible probabilistic networks, despite a computationally expensive parameter estimation. In Olde Loohuis *et al.* (2014), we introduced an algorithm called *CAnceR PRogression Extraction with Single Edges* (CAPRESE), which, based on its extensive empirical analysis, may be deemed as the current state-of-the-art algorithm for the inference of tree models of cancer progression. It is based on a shrinkage-like statistical estimation, grounded in a general theoretical framework, which we extend further in this paper. Other results that extend tree representations of cancer evolution exploit mixture tree models, i.e. multiple oncogenetic trees, each of which can independently result in cancer development (Beerenwinkel *et al.*, 2005). In general, all these methods are capable of modeling diverging temporal orderings of events in terms of branches, although the possibility of converging evolutionary paths is precluded.

To overcome this limitation, the most recent approaches tend to adopt Bayesian graphical models, i.e. Bayesian Networks (BN). In the literature, there have been two initial families of methods aimed at inferring the structure of a BN from data (Koller and Friedman, 2009). The first class of models seeks to explicitly capture all the conditional independence relations encoded in the edges and will be referred to as *structural approaches*; the methods in this family are inspired by the work on causal theories by Judea Pearl (cf., Pearl, 1988, 2000; Spirtes *et al.*, 2000; Tsamardinos *et al.*, 2003).

The second class—*likelihood approaches*—seeks a model that maximizes the likelihood of the data (cf., Carvalho, 2009; Heckerman *et al.*, 1995; Schwarz, 1978).

A more recent *hybrid approach* to learn a BN which combines the two families above by (i) constraining the search space of the valid solutions and, then, (ii) fitting the model with likelihood maximization (see Beerenwinkel *et al.*, 2007; Gerstung *et al.*, 2009; Misra *et al.*, 2014). A further technique to reconstruct progression models from cross-sectional data was introduced in Attolini *et al.* (2010), in which the transition probabilities between genotypes are inferred by defining a Moran process that describes the evolutionary dynamics of mutation accumulation. In Cheng *et al.* (2012) this methodology was extended to account for pathway-based phenotypic alterations.

## 2 Approach

In what follows, we denote with  $P(\cdot)$  and  $P(\cdot|\cdot)$  the observed marginal and conditional probability of an event, whose complement is denoted with the diacritical mark  $\bar{\cdot}$  (macron).

*A probabilistic model of selective advantage.* Central to CAPRI's score function is Suppes' notion of *probabilistic causation* (Suppes, 1970), which can be stated in the following terms: a selectivity relation holds among two observables *i* and *j* if (i) *i* occurs earlier than *j*—*temporal priority* (TP)—and (ii) if the probability of observing *i* raises the probability of observing *j*, i.e.  $P(j|i) > P(j|\bar{i})$ —*probability raising* (PR). (Suppes presents the relation in terms of causality; however, we avoid Suppes' terminology as we build on just two of his many axioms, which only give rise to the notion of *prima facie* causality.) The definition of probability raising subsumes positive statistical dependency and mutuality (see Olde Loohuis *et al.*, 2014). Note that the resulting relation (also called *prima facie* causality) is purely

observational and remains agnostic to the possible mechanistic cause-effect relation involving  $i$  and  $j$ .

While Suppes' definition of probabilistic causation has known limitations in the context of general causality theory (see discussions in, e.g. Hitchcock, 2012; Kleinberg, 2012), in the context of cancer evolution, this relation appropriately describes various features of *selective advantage* in somatic alterations that accumulate as tumor progresses.

Thus, in our framework, we implement the temporal priority among events—condition (1)—as  $\mathcal{P}(i) > \mathcal{P}(j)$ , because it is intuitively sound to assume that the (cumulative) genomic events occurring earlier are the ones present in higher frequency in a dataset. In addition, condition (2) is implemented as is, that is by requiring that for each pair of observables  $i$  and  $j$  directly connected,  $\mathcal{P}(j|i) > \mathcal{P}(j|\bar{i})$  is verified. Taken together, these conditions give rise to a natural ordering relation among events, written ' $i \triangleright j$ ' and read as ' $i$  has a selective influence on  $j$ .' This relation is a *necessary* but *not sufficient* condition to capture the notion of selective advantage, and additional constraints need to be imposed to filter spurious relations. Spurious correlations are both intrinsic to the definition (e.g. if  $i \triangleright j \triangleright w$  then also  $i \triangleright w$ , which could be spurious) and to the model we aim at inferring, because data is finite as well as corrupted by noise.

Building on this framework, we devise inference algorithms that capture the essential aspects of heterogeneous cancer progressions: *branching*, *independence* and *convergence*—all combining in a progression model.

**Progression patterns.** The complexity of cancer requires modeling multiple non-trivial *patterns* of its progression: for a specific event, a pattern is defined as a specific combination of the closest upstream events that confers a selective advantage.

As an example, imagine a clonal subpopulation becoming fit—thus enjoying expansion and selection—once it acquires a mutation of gene  $c$ , provided it also has previously acquired a mutation in a gene in the upstream  $alb$  pathway. In terms of progression, we would like to capture the trajectories:  $\{a, \neg b\}$ ,  $\{\neg a, b\}$  and  $\{a, b\}$  precedes  $c$  (where  $\neg$  denotes the absence of an event in the gene).

To establish this analysis formally, we augment our model of selection in a tumor with a language built from simple propositional logic formulas using the usual Boolean connectives: namely, 'and' ( $\wedge$ ), 'or' ( $\vee$ ) and 'xor' ( $\oplus$ ). These patterns can be described by formulae in a propositional logical language, which can be rendered in *conjunctive normal form* (CNF). A CNF formula  $\phi$  has the following syntax:  $\phi = c_1 \wedge \dots \wedge c_n$ , where each  $c_i$  is a *disjunctive clause*  $c_i = c_{i,1} \vee \dots \vee c_{i,k}$  over a set of literals, each literal representing an event or its negation. Given this (rather obvious) pattern representation, we write the conditions for *selectivity with patterns* as

$$\phi \triangleright e \Leftrightarrow \mathcal{P}(\phi) > \mathcal{P}(e) \text{ and } [\mathcal{P}(e|\phi) > \mathcal{P}(e|\bar{\phi})]; \quad (1)$$

with respect to the example above, patterns could be  $a \vee b \triangleright c$  and  $a \oplus b \triangleright c$ . (Note that the conjunction  $\wedge$  in our setting is interpreted differently from the classical notion [and the one adopted in e.g. Gerstung et al., 2009], since  $a \wedge b \triangleright c$  implies  $a \triangleright c$  and  $b \triangleright c$  in our framework. See also Beerenwinkel et al., 2014. Moreover, note that the scope of this study is intentionally kept limited from further generalization of formulae i.e. we will not consider statements of the form  $\phi_i \triangleright \phi_j$ , where the rightmost argument is a formula too.)

In our framework the problem of reconstructing a probabilistic graphical model of progression reduces to the following: for each input event  $e$ , assess a *set of selectivity patterns*  $\{\phi_1 \triangleright e, \dots, \phi_k \triangleright e\}$ , filter the spurious ones, and combine the rest

in a *direct acyclic graph* (DAG), augmented with logical symbols. (A DAG is formed by a set of nodes and oriented edges connecting one node to another, such that there are no directed loops among them. See [Supplementary information, Section 1](#) for a technical definition.) Notice that while we broke down the progression extraction into a series of sub-tasks, the problem remains complex: patterns are unknown, potentially spurious and exponential in formula size; data is noisy; patterns must allow for 'imperfect regularities', rather than being strict. (This statement implies that there could be samples—i.e. patients—contradicting a pattern which still remains valid at a population level. For this reason a pattern  $x \wedge y \triangleright z$  is sometimes called a 'noisy and'.) To summarize, in our setting we can model complex progression trajectories with branches (i.e. events involved in various patterns), independent progressions (i.e. events without common ancestors) and convergence (via CNF formulas). The framework we introduce here is highly versatile, and to the best of our knowledge, it infers and checks more complex claims than any cancer progression algorithms described thus far (cf., Desper et al., 1999; Gerstung et al., 2009; Olde Loohuis et al., 2014).

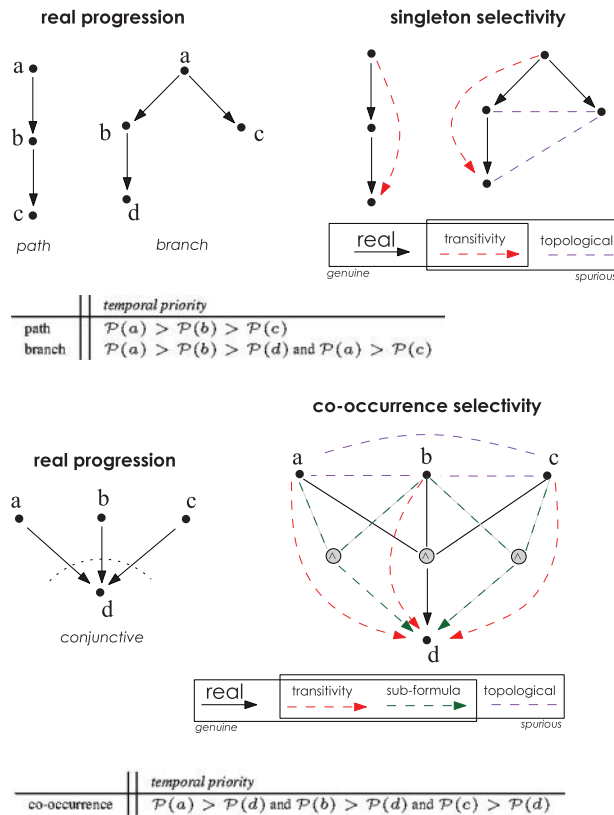
### 3 Methods

Building on the framework described in the previous section, we now describe the implementation of CAPRI's building blocks comparing it against other state-of-the-art approaches, as well as through an application involving reconstruction of a leukemia progression model. Notice that, in general, the inference of cancer progression models requires a complex *data processing pipeline*, as summarized in [Figure 3](#); its architecture optimally exploits CAPRI's efficiency.

**Assumptions.** CAPRI relies on the following assumptions: (i) Every pattern is expressible as a propositional CNF formula; (ii) All events are persistent, i.e. an acquired mutation cannot disappear; (iii) All relevant events in tumor progression are observable, with the observations describing the progressive phenomenon in an essential manner (i.e. *closed world* assumption, in which all events 'driving' the progression are detectable); (iv) All the events have non-degenerate observed probability in  $(0, 1)$ ; (v) All events are distinguishable, simultaneously observed nor simultaneously unobservable in the following sense: input alterations produce different profiles across input samples. Assumptions (i–ii) relate to the framework derived in previous section, while (iii) imposes an onerous burden on the experimentalists, who must select the *relevant* genomic events to model. (Theoretically, this assumption—common to other Bayesian learning problems—is necessary to prove CAPRI's ability to extract the exact model in the optimal case of infinite samples. Practically, as all *relevant* events are hardly selectable a priori and sample size is finite, further statistics can be used to select the most relevant driver alterations—see also Section 4, Results and Discussion. Nonetheless, CAPRI can provide significant results even if this assumption is not or cannot be verified.) Assumption (iv) relates instead to the statistical distinguishability of the input events (see the next section on CAPRI's Data Input).

**Trading Complexity for Expressivity.** To automatically extract the patterns that underlie a progression model, one may try to adopt a brute-force method of enumerating and testing all possibilities. This strategy is computationally intractable, however, since the number of (distinct) (sub)formulae grows exponentially with the number of events included in the model. Therefore, we need to exploit certain





**Fig. 2.** Singleton and co-occurrence selectivity patterns. Examples of patterns that CAPRI can automatically extract without prior hypotheses. (Top): A linear path and branching model (left) and corresponding singleton selectivity patterns with infinite sample size (right). All the genuine connections are shown (red and black, directed by the temporal priority), as well as edges (purple, undirected) which might be suggested by the topology (or observations, if data were finite). (Bottom): Example of conjunctive model ( $a$  and  $b$  and  $c$ ). The co-occurrence selectivity pattern is shown, with all true patterns and infinite sample size. The topology is augmented by logical connectives; green arrows are spurious patterns emerging from the structure of the true pattern  $a \wedge b \wedge c \triangleright d$

properties of the  $\triangleright$  relation whenever possible, and trade expressivity for complexity in other cases, as explained below.

Note that *singleton* and *co-occurrence* ( $\wedge$ ) types of patterns are amenable to *compositional reasoning*: if  $i_1 \wedge \dots \wedge i_k \triangleright j$  then, for any  $p = 1, \dots, k$ ,  $i_p \triangleright j$ . This observation leads to the following straightforward strategy of evaluating every conjunctive (and henceforth singleton) relation using a pairwise-test for the selectivity relation (see Fig. 2).

Unfortunately, it is easy to see that this reasoning fails to generalize for CNF patterns: e.g. when the pattern contains disjunctive operators ( $\vee$ ). As an example, consider pattern  $a \vee b \triangleright c$ , in a cancer where  $\{a, \neg b\}$  progression to  $c$  is more prevalent than  $\{\neg a, b\}$  and  $\{a, b\}$ . In this case, considering sub-formulas only we might find  $a \triangleright c$  but miss  $b \triangleright c$  because the probability of mutated  $b$  is smaller than that of  $c$ , thus invalidating condition (1) of relation  $\triangleright$ . Notice that in extreme situations, when the data is very noisy, the algorithm may even ‘invert’ the selectivity relation to  $c \triangleright b$ .

This difficulty is not a peculiarity of our framework, but rather intrinsic to the problem of extracting complex ‘causal networks’ (cf., Pearl, 1988, 2000; Kleinberg, 2012). To handle this situation, CAPRI adapts a strategy that trades complexity for expressivity: Figure 3 and described as the resulting inference procedure,

Algorithm 1, can be executed in two modes: unsupervised and supervised. In the former, inferred patterns of confluent progressions are constrained to co-occurrence types of relations, in the latter CAPRI can test more complex patterns, i.e. disjunctive or ‘mutual exclusive’ ones, provided they are given as prior hypotheses. In both cases, CAPRI’s complexity—studied in next sections—is quadratic both in the number of events and hypotheses.

**Data Input (Step 1).** CAPRI (cf., Algorithm 1) requires an input set  $G$  of  $n$  events, i.e. genomic alterations, and  $m$  cross-sectional samples, represented as a dataset in an  $m \times n$  binary matrix  $D$ , in which an entry  $D_{i,j} = 1$  if the event  $j$  was observed in sample  $i$ , and 0 otherwise. Assumption (iv) is satisfied when all columns in  $D$  differ—i.e. the alteration profiles yield different observations.

Optionally, a set of  $k$  input hypotheses  $\Phi = \{\varphi_1 \triangleright e_1, \dots, \varphi_k \triangleright e_k\}$ , where each  $\varphi_i$  is a well-formed CNF formula. (Formally, we require that  $\varphi_i \sqsubseteq e_i$ , where  $\sqsubseteq$  represents the usual *syntactical* ordering relation among atomic events and formulas, and disallows for example  $a \vee b \triangleright a$ .) Note that we advise that the algorithm be used in the following regime:  $k + n \gg m$ . (In the current biomedical setting, the number of samples ( $m$ ) is usually in the hundreds, while number of possible mutations ( $n$ ) and hypotheses ( $k$ ), absent any pre-processing, could be large, thus violating the assumption; in these cases, we rely on various commonly used pre-preprocessing filters to limit  $n$  to driver mutations, and  $k$  to simple hypotheses involving the driver mutations. However, in the future as the number of samples increases, we envision a more agnostic application.)

**Data Preprocessing (Lifting, step 2).** When input hypotheses are provided (e.g. by a domain expert), CAPRI first performs a *lifting operation* over  $D$  to permit direct inference of complex selectivity relations over a joint representation, which involve input events as well as the hypotheses. Lifting operation evaluates each input CNF formula—for all input hypotheses in  $\Phi$ —and outputs a lifted matrix  $D(\Phi)$  to be processed further as in step 1. As an example, consider hypothesis  $a \oplus b \triangleright c$  lifted input matrix  $D$  is:

$$D(\Phi) = \begin{bmatrix} a & b & c & a \oplus b \triangleright c \\ 1 & 1 & 1 & 1 \oplus 1 = 0 \\ 1 & 0 & 1 & 1 \oplus 0 = 1 \\ 0 & 1 & 0 & 0 \oplus 1 = 1 \\ 1 & 0 & 1 & 1 \oplus 0 = 1 \end{bmatrix}.$$

Note that the first row (profile  $\{a, b, c\}$ ) contradicts the hypothesis, while all other rows support it.

**Selectivity Topology (steps 3–5).** We exploit a compositional approach to test CNF hypotheses as follows: the disjunctive relations are grouped, and treated as if they were individual objects in  $G$ . For example, when a formula  $\varphi \triangleright d$  where  $\varphi = (a \vee b) \wedge c$  is considered, we assess  $\varphi \triangleright d$  as whether  $(a \vee b) \triangleright d$  and  $c \triangleright d$  hold—with the proviso that we treat  $(a \vee b)$  as an individual event. Formally, with clauses( $\varphi$ ) we denote the disjunctive clauses in a CNF formula.

Nodes in the reconstruction are all input events together with all the disjunctive clauses of each input formula  $\varphi$ .

Edges in the reconstructed DAG are patterns that satisfy both conditions (1) and (2) of the selectivity relation  $\triangleright$ . Formally, CAPRI includes an edge between two nodes  $\varphi$  and  $j$  only if both  $\Gamma_{\varphi,j} = \mathcal{P}(\varphi) - \mathcal{P}(j)$  and  $\Lambda_{\varphi,j} = \mathcal{P}(j|\varphi) - \mathcal{P}(j|\neg\varphi)$  are strictly positive. Note that  $\varphi$  can be both a disjunctive clause as well as a singleton event. A function  $\pi(\cdot)$  assigns a parent to each node that is not an input formula. Note that this approach works efficiently by nature

of the lifted representation of  $D$ . The reconstructed DAG contains all the true positive patterns, with respect to  $\triangleright$ , plus spurious instances of  $\triangleright$  which CAPRI subsequently removes in step 6 (cf., the [Supplementary Material](#) for a proof of this statement).

Note that  $\mathcal{D}$  can be readily interpreted as a probabilistic graphical model, once it is augmented with a labeling function  $\alpha : N \rightarrow [0, 1]$ , where  $N$  is the set of nodes—i.e. the genetic alterations—such that  $\alpha(i)$  is the *independent probability* of observing mutation  $i$  in a sample, whenever *all of its parent* mutations (i.e.  $\pi(i)$ ) are observed (if any). Thus  $\mathcal{D}$  induces a *distribution* of observing a subset of events in a set of samples (i.e. a probability of observing a certain *mutational profile* in a patient).

**Maximum Likelihood Fit (step 6).** As the selectivity relation provides only a *necessary* condition, we must filter out all of its *spurious instances* that might have been included in  $\mathcal{D}$  (i.e. the possible *false positives*).

For any selectivity structure, spurious claims contribute to a reduction in the *likelihood-fit* relative to true patterns. Thus, a standard maximum-likelihood fit can be used to select and prune the selectivity DAG (including a *regularization term* to avoid overfitting. [In principle other regularization strategies common to Bayesian learning could be used, e.g. Akaike information criterion (see [Carvalho, 2009](#) and references therein). In this article, we prefer to work with BIC which, in general, trades model complexity to reduce false positives rate.]). Here, we adopt the *Bayesian Information Criterion* (BIC), which implements *Occam's razor* by combining log-likelihood fit with a *penalty criterion* proportional to the log of the DAG size via *Schwarz Information Criterion* (see [Schwarz, 1978](#)). The BIC score is defined as follows.

$$\text{bic}(\mathcal{D}, D(\Phi)) = \mathcal{LL}(\mathcal{D}, D(\Phi)) - \frac{\log m}{2} \dim(\mathcal{D}). \quad (2)$$

Here,  $D(\Phi)$  is the lifted input matrix,  $m$  denotes the number of samples and  $\dim(\mathcal{D})$  is the number of parameters in the model  $\mathcal{D}$ . Because, in general,  $\dim(\cdot)$  depends on the number of parents each node has, it is a good metric for model complexity. Moreover, since each edge added to  $\mathcal{D}$  increases model complexity, the regularization term based on  $\dim(\cdot)$  favors graphs with fewer edges and, more specifically, fewer parents for each node.

At the end of this step,  $\mathcal{D}$  and the labeling function are modified accordingly, based on the result of BIC regularization. By collecting all the incoming edges in a node it is possible to extract the patterns, which have been selected by CAPRI as the positive ones.

---

#### Algorithm 1 Cancer Progression InfERENCE (CAPRI)

- 1: **Input:** A set of events  $G = \{g_1, \dots, g_n\}$ , a matrix  $D \in \{0, 1\}^{m \times n}$  and  $k$  CNF causal claims  $\Phi = \{\varphi_1 \triangleright e_1, \dots, \varphi_k \triangleright e_k\}$  where, for any  $i$ ,  $e_i \sqsubseteq \varphi_i$  and  $e_i \in G$ ;
- 2: [*Lifting*] Define the *lifting* of  $D$  to  $D(\Phi)$  as the augmented matrix

$$D(\Phi) = \left[ \begin{array}{ccc|ccc} D_{1,1} & \dots & D_{1,n} & \varphi_1(D_{1,\cdot}) & \dots & \varphi_k(D_{1,\cdot}) \\ \vdots & \ddots & \vdots & \vdots & \ddots & \vdots \\ D_{m,1} & \dots & D_{m,n} & \varphi_1(D_{m,\cdot}) & \dots & \varphi_k(D_{m,\cdot}) \end{array} \right].$$

by adding a column for each  $\varphi_i \triangleright e_i \in \Phi$ , with  $\varphi_i$  evaluated row-by-row. Define then the coefficients  $\Gamma_{i,j} = \mathcal{P}(i) - \mathcal{P}(j)$  and  $\Lambda_{i,j} = \mathcal{P}(j|i) - \mathcal{P}(j|\bar{i})$  pairwise over  $D(\Phi)$ ;

- 3: [*DAG nodes*] Define the set of nodes  $N = G \cup (\cup_{\varphi_i} \text{clauses}(\varphi_i))$  which contains both input events and the disjunctive clauses in every input formula of  $\Phi$ .
- 4: [*DAG edges*] Define a parent function  $\pi$  where  $\pi(j \in G) = \emptyset$ —avoid edges incoming in a formula and

$$\pi(j \in G) = \{i \in G \mid \Gamma_{i,j}, \Lambda_{i,j} > 0\} \cup \{\text{clauses}(\varphi) \mid \Gamma_{\varphi,j}, \Lambda_{\varphi,j} > 0, \varphi \triangleright j \in \Phi\}. \quad (3)$$

Set the DAG to  $\mathcal{D} = (N, \pi)$ . (Although CAPRI is equipped with bootstrap testing it is still possible to encounter various degenerate situations. In particular, for some pair of events it could be that temporal priority cannot be satisfactorily resolved, i.e. there is no significant  $P$ -value for any edge orientation. Thus, loops might be present in the inferred *prima facie* topology. Nonetheless, some of these could be still disentangled by probability raising, while some might remain, albeit rarely. To remove such edges we suggest to proceed as follows: (i) sort these edges according to their  $P$ -value (considering both temporal priority and probability raising), (ii) scan the sorted list in decreasing order of confidence, (iii) remove an edge if it forms a loop.)

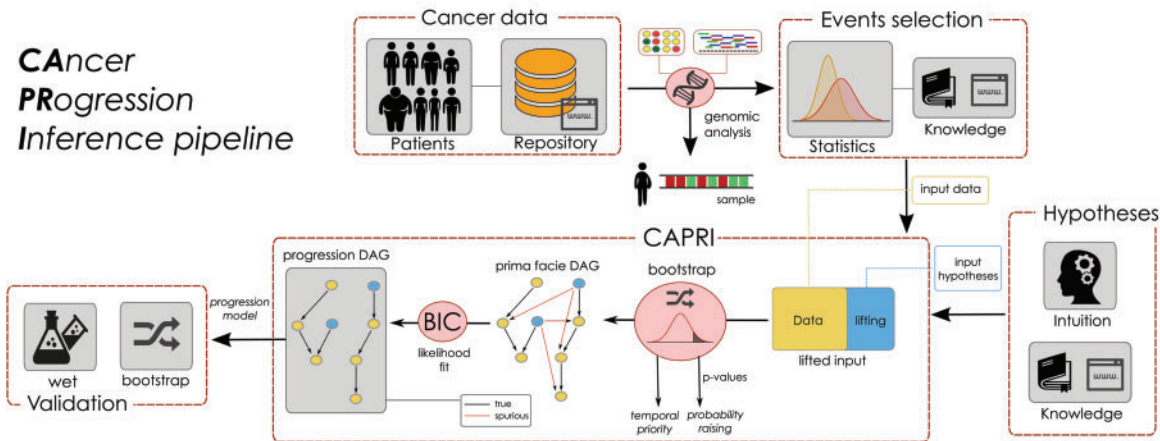
- 5: [*DAG labeling*] Define the labeling  $\alpha$  as follows

$$\alpha(j) = \begin{cases} \mathcal{P}(j), & \text{if } \pi(j) = \emptyset \text{ and } j \in G; \\ \mathcal{P}(j|i_1 \wedge \dots \wedge i_n), & \text{if } \pi(j) = \{i_1, \dots, i_n\}. \end{cases}$$

- 6: [*Likelihood fit*] Filter out all spurious causes from  $\mathcal{D}$  by likelihood fit with the regularization BIC score and set  $\alpha(j) = 0$  for each removed edge.
  - 7: **Output:** the DAG  $\mathcal{D}$  and  $\alpha$ ;
- 

**Inference Confidence: Bootstrap and Statistical Testing.** To infer confidence intervals of the selectivity relations  $\triangleright$ , CAPRI employs *bootstrap with rejection resampling* as follows, by estimating a distribution of the marginal and joint probabilities. For each event, (i) CAPRI samples with repetitions rows from the input matrix  $D$  (bootstrapped dataset), (ii) CAPRI next estimates the distributions from the observed probabilities, and finally, (iii) CAPRI rejects values which do not satisfy  $0 < \mathcal{P}(i) < 1$  and  $\mathcal{P}(i|j) < 1 \vee \mathcal{P}(j|i) < 1$ , and iterates restarting from (i). We stop when we have, for each distribution, at least  $K$  values (in our case  $K = 100$ ). Any inequality (i.e. checking temporal priority and probability raising) is estimated using the non-parametric Mann-Whitney U test with  $P$ -values set to 0.05. (The Mann-Whitney U test is a rank-based non-parametric statistical hypothesis test that can be used as an alternative to the Student's t-test and is particularly useful if data are not normally distributed.) We compute confidence  $P$ -values for both temporal priority and probability raising using this test, which need not assume Gaussian distributions for the populations.

Once a DAG  $\mathcal{D}$  is inferred both *parametric* and *non-parametric bootstrapping methods* can be used to assign a confidence level to its respective pattern and to the overall model. Essentially, these tests consist of using the reconstructed model (in the parametric case), or the probabilities observed in the dataset (in the non-parametric case) to generate new synthetic datasets, which are then reused to reconstruct the progressions (see, e.g. [Efron, 2010](#) for an overview of these methods). The confidence is estimated by the number of times the DAG or any instance of  $\triangleright$  is reconstructed from the generated data.



**Fig. 3.** Data processing pipeline for cancer progression inference. We sketch a pipeline to best exploit CAPRI's ability to extract cancer progression models from cross-sectional data. Initially, one collects *experimental data* (which could be accessible through publicly available repositories such as TCGA) and performs *genomic analyses* to derive profiles of, e.g. somatic mutations or Copy-Number Variations for each patient. Then, statistical analysis and biological priors are used to select events relevant to the progression and imputable by CAPRI—e.g. *driver mutations*. To exploit CAPRI's supervised execution mode (see Methods) one can use further statistics and priors to generate *patterns of selective advantage*—e.g. hypotheses of mutual exclusivity. CAPRI can extract a progression model from these data and assess various *confidence* measures on its constituting relations—e.g. (non-)parametric bootstrap and hypergeometric testing. *Experimental validation* concludes the pipeline

*Complexity, Correctness and Expressivity.* CAPRI has the following asymptotic complexity (Theorem 1, [Supplementary information Section 2](#)):

- Without input hypotheses the execution is self-contained and polynomial in the size of  $D$ .
- In addition to the above cost, CAPRI tests input hypotheses of  $\Phi$  at a polynomial cost in the size of  $|\Phi|$ . In this case, however, its complexity may range over many orders of magnitude depending on the structural complexity of the input set  $\Phi$  consisting of hypotheses.

An empirical analysis of the execution time of CAPRI and the competing techniques on synthetic datasets is provided in the [Supplementary information, Section 3.5](#).

CAPRI is a *sound and complete* algorithm, and its expressivity in terms of the inferred patterns is proportional to the hypothesis set  $\Phi$  which, in turn, determines the complexity of the algorithm. With a proper set of input hypothesis, CAPRI can infer all (and only) the true patterns from the data, filtering out all the spurious ones (Theorem 2, [Supplementary information Section 2](#)). Without hypotheses, besides singleton and co-occurrence, no other patterns can be inferred (see [Fig. 2](#)). Also, some of these claims might be spurious in general for more complex (and unverified) CNF formula (Theorem 3, [Supplementary information Section 2](#)).

## 4 Results and discussion

To determine CAPRI's relative accuracy (true-positives and false-negatives) and performance compared with the state-of-the-art techniques for *network inference*, we performed extensive *simulation experiments*. From a list of potential competitors of CAPRI, we selected: *Incremental Association Markov Blanket* (IAMB, [Tsamardinos et al., 2003](#)), the *PC algorithm* (see [Spirtes et al., 2000](#)), *Bayesian Information Criterion* (BIC, [Schwarz, 1978](#)), *Bayesian Dirichlet with likelihood equivalence* (BDE, [Heckerman et al., 1995](#)), *Conjunctive Bayesian Networks* (CBN, [Gerstung et al., 2009](#)) and *Cancer Progression Inference with Single Edges*

(CAPRESE, [Olde Loohuis et al., 2014](#)). These algorithms constitute a rich landscape of structural methods (IAMB and PC), likelihood scores (BIC and BDE) and hybrid approaches (CBN and CAPRESE).

Also, we applied CAPRI to the analysis of an atypical Chronic Myeloid Leukemia dataset of somatic mutations with data based on [Piazza et al. \(2013\)](#).

### 4.1 Synthetic data

We performed extensive tests on a large number of *synthetic datasets* generated by randomly parametrized progression models with distinct key features, such as the presence/absence of: (i) *branches*, (ii) *confluences with patterns of co-occurrence*, (iii) *independent progressions* (i.e. composed of disjoint sub-models involving distinct sets of events). Accordingly, we distinguish four classes of generative models with increasing complexity and the following features:

	trees	Forests	connected DAGs	disconnected DAGs
(1)	✓	✓	✓	✓
(2)	✗	✗	✓	✓
(3)	✗	✓	✗	✓

The choice of these different type of topologies is not a mere technical exercise, but rather it is motivated, in our application of primary interest, by *heterogeneity of cancer cell types* and *possibility of multiple cells of origin*.

To account for *biological noise* and *experimental errors* in the data we introduce a parameter  $\nu \in (0, 1)$  which represents the probability of each entry to be random in  $D$ , thus representing a *false positive* ( $\epsilon_+$ ) and a *false negative* rate ( $\epsilon_-$ ):  $\epsilon_+ = \epsilon_- = \nu/2$ . The noise level complicates the inference problem, since samples generated from such topologies will likely contain sets of mutations that are correlated but causally irrelevant.

To have reliable statistics in all the tests, 100 distinct progression models per topology are generated and, for each model, for every chosen combination of sample set size  $m$  and noise rate  $\nu$ , 10 different datasets are sampled (see [Supplementary information Section 3](#) for our synthetic data generation methods).

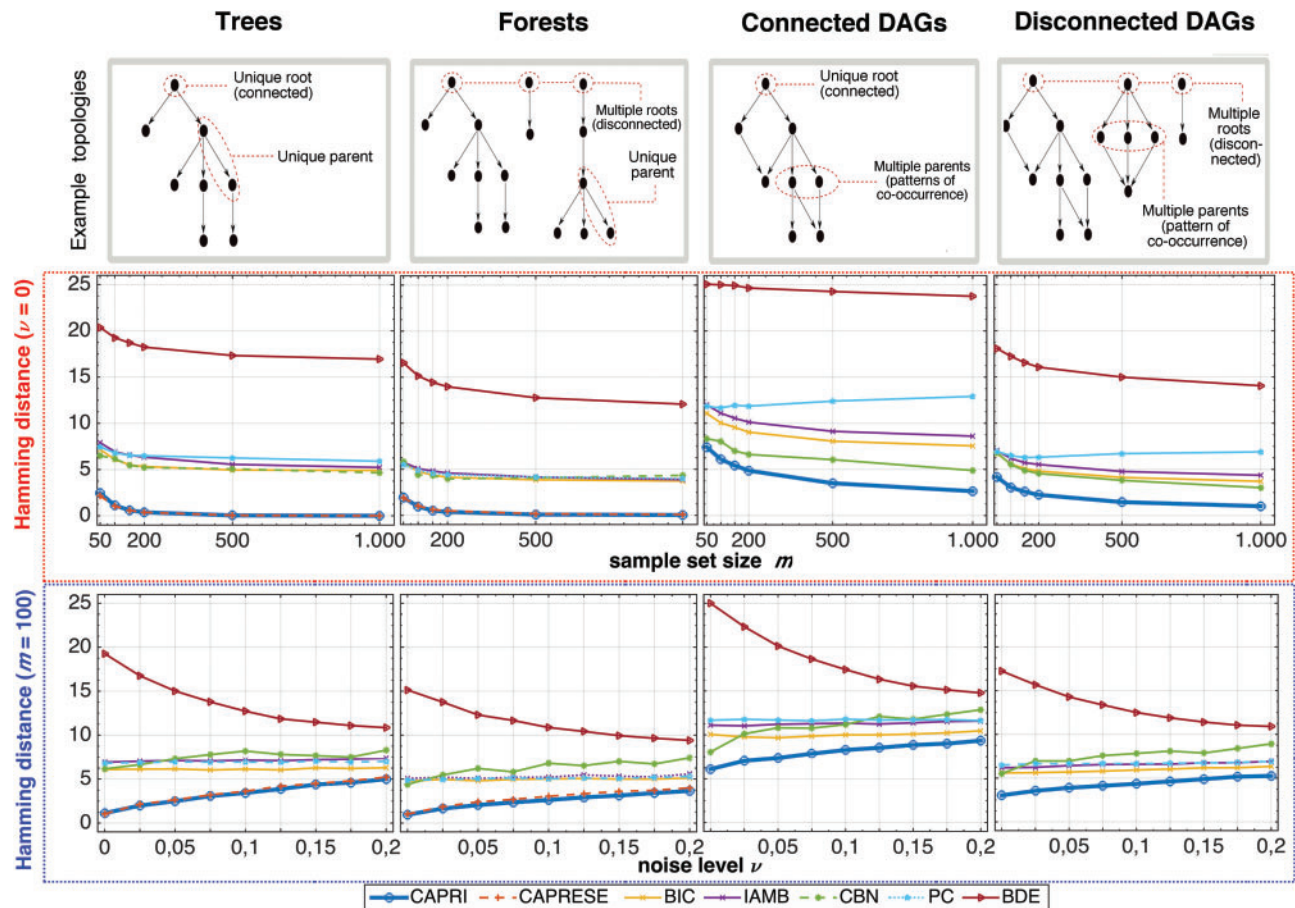


Algorithmic performance was evaluated using the metrics *Hamming distance* (HD), *precision* and *recall*, as a function of data-set size,  $\epsilon_+$  and  $\epsilon_-$ . HD measures the *structural similarity* among the reconstructed progression and the generative model in terms of the minimum-cost sequence of node edit operations (inclusion and exclusion) that transforms the reconstructed topology into the generative one. (This measure corresponds to the sum of false positives and false negative and, for a set of  $n$  events, is bounded above by  $n(n-1)$  when the reconstructed topology contains all the false negatives and positives.) Precision and recall are defined as follows:  $\text{precision} = \text{TP}/(\text{TP} + \text{FP})$  and  $\text{recall} = \text{TP}/(\text{TP} + \text{FN})$ , where TP are the *true positives* (number of correctly inferred true patterns), FP are the *false positives* (number of spurious patterns inferred) and FN are the *false negatives* (number of true patterns that are *not* inferred). The closer both precision and recall are to 1, the better.

In Figure 4 we show the performance of CAPRI and of the competing techniques, in terms of Hamming distance, on datasets generated from models with 10 events and all the four different topologies. In particular, we show the performance: (i) in the case of noise-free datasets, i.e.  $\nu = 0$  and different values of the sample set size  $m$  and (ii) in the case of a fixed sample set size,  $m = 100$  (size

that is likely to be found in currently available cancer databases, such as TCGA (cf., NCI and the NHGRI, 2005)) and different values of the noise rate  $\nu$ . As is evident from Figure 4 CAPRI outperforms all the competing techniques with respect to all the topologies and all the possible combinations of noise rate and sample set size, in terms of average Hamming distance (with the only exception of CAPRESE in the case of tree and forests, which displays a behavior closer to CAPRI's). The analyses on precision and recall display consistent results (Supplementary information Section 3). In other words, we demonstrate on the basis of extensive synthetic tests that CAPRI requires a much lower number of samples than the other techniques to converge to the real generative model and also that it is much more robust even in the presence of significant amount of noise in the data, irrespective of the underlying topology.

See Supplementary information Section 3 for a more complete description of the performance evaluation for all the analyzed combinations of parameters. There, we have shown that CAPRI is highly effective when the co-occurrence constraint on confluent is relaxed to *disjunctive* patterns, *even if no input hypotheses are provided*, i.e.  $\Phi = \emptyset$ . This result hints at CAPRI's robustness to infer patterns with imperfect regularities. Finally, we also show that



**Fig. 4.** Comparative study. Performance and accuracy of CAPRI (unsupervised execution) and other algorithms, IAMB, PC, BIC, BDE, CBN and CAPRESE, were compared using synthetic datasets sampled by a large number of randomly parametrized progression models—*trees*, *forests*, *connected* and *disconnected DAGs*, which capture different aspects of confluent, branched and heterogenous cancer progressions. For each of those, 100 models with  $n = 10$  events were created and 10 distinct datasets were sampled by each model. Datasets vary by number of samples ( $m \in \{50, 100, 150, 200, 500, 1000\}$ ), when data contain no noise ( $\nu = 0$ ). The lower the HD, the smaller is the total rate of mis-inferred selectivity relations among events. (Red box) Average Hamming distance (HD)—with 1000 runs—between the reconstructed and the generative model, as a function of dataset size ( $m \in \{50, 100, 150, 200, 500, 1000\}$ ), when data contain no noise ( $\nu = 0$ ). The lower the HD, the smaller is the total rate of mis-inferred selectivity relations among events. (Blue box) The same is shown for a fixed sample set size  $m = 100$  as a function of noise level in the data ( $\nu \in \{0, 0.025, 0.05, \dots, 0.2\}$ ) so as to account for input *false positives* and *false negatives*. See Supplementary information Section 3 for more extensive results on precision and recall scores and also including additional combinations of noise and samples as well as experimental settings



CAPRI is effective in inferring synthetic lethality relations in this case using the operator  $\oplus$  as introduced in Section 2, Approach; when a combination of mutations in two or more genes leads to cell death, while separately, the mutations are viable. In this case, candidate relations are directly input as  $\Phi$ .

## 4.2 Atypical Chronic Myeloid Leukemia

As a case study, we applied CAPRI to the mutational profiles of 64 ACML patients described in Piazza et al. (2013). Through exome sequencing, the authors identify a recurring *missense point mutation* in the *SET-binding protein 1* (SETBP1) gene as a novel ACML marker.

Among all the genes present in the dataset by Piazza et al., we selected those either (i) mutated—considered any mutation type—in at least 5% of the input samples (3 patients) or (ii) hypothesised to be part of a functional ACML progression pattern in the literature. (Two *hard exclusivity* patterns—i.e. mutual exclusivity with ‘xor’—were tested, involving the mutations of: (i) genes ASXL1 and SF3B1 [see Lin et al., 2014], which is present in the inferred progression model in Fig. 5 and (ii) genes TET2 and IDH2 [see Figueroa et al., 2010]. The syntax expressing the patterns is as described in the Supplementary information, Section 4.) The input dataset with selected events is shown in Figure 5; notice that somatic mutations are categorised as *indel*, *missense point* and *nonsense point* as in Piazza et al. (2013). In Figure 5 we show the model reconstructed by CAPRI (supervised mode, execution time  $\approx 5$  seconds) on this dataset, with confidence estimated via 1000 non-parametric bootstrap iterations. The model highlights several non trivial selectivity relations involving genomic events relevant to ACML development.

First, CAPRI predicts a progression involving mutations in SETBP1, ASXL1 and CBL, consistently with the recent study by Meggendorfer et al. (2013), in which these genes were shown to be highly correlated and possibly functioning in a synergistic manner for ACML progression. Specifically, CAPRI predicts a selective advantage relation between missense point mutations in SETBP1 and nonsense point mutations in ASXL1. This is in line with recent evidence from Inoue et al. (2014) suggesting that SETBP1 mutations are enriched among ASXL1-mutated *myelodysplastic syndrome* (MDS) patients, and *in-vivo* experiments point to a driver role of SETBP1 for that leukemic progression. Interestingly, our model seems also to suggest a different role of ASXL1 *missense* and *nonsense* mutation

types in the progression, yet more extensive studies (e.g. prospective or systems biology explanation) are needed to corroborate this hypothesis.

Among the hypotheses given as input to CAPRI, the algorithm seems to suggest that the exclusivity pattern among ASXL1 and SF3B1 mutations selects for CBL missense point mutations. The role of the ASXL1/SF3B1 exclusivity pattern is consistent with the study of Lin et al. (2014) which shows that, on a cohort of 479 MDS patients, mutations in SF3B1 are inversely related to ASXL1 mutations.

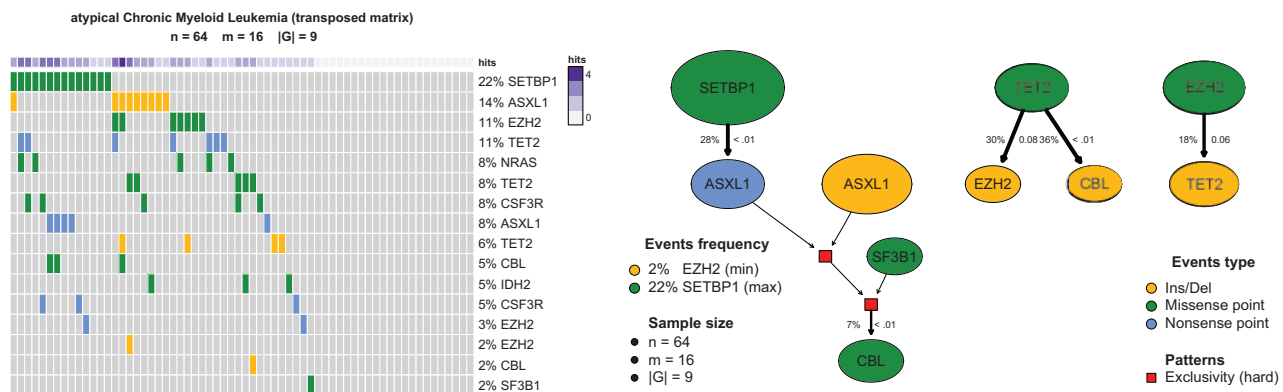
Also, in Abdel-Wahab et al. (2012) it was recently shown that ASXL1 mutations, in patients with MDS, *myeloproliferative neoplasms* (MPN) and *acute myeloid leukemia*, most commonly occur as nonsense and insertion/deletion in a clustered region adjacent to the highly conserved PHD domain (see Gelsi-Boyer et al., 2009) and that mutations of any type eventually result in a loss of ASXL1 expression. This observation is consistent with the exclusivity pattern among ASXL1 mutations in the reconstructed model, possibly suggesting alternative trajectories of somatic evolution for ACML (involving either ASXL1 nonsense or indel mutations).

Finally, CAPRI predicts selective advantage relations among TET2 and EZH2 missense point and indel mutations. Even though the limited sample size does not permit drawing definitive conclusions on the ordering of such alterations, we can hypothesize that they may play a synergistic role in ACML progression. Indeed, Muto et al. (2013) suggests that the concurrent loss of EZH2 and TET2 might cooperate in the pathogenesis of myelodysplastic disorders, by accelerating the overall tumor development, with respect to both MDSs and *overlap disorders* (MDS/MPN).

## 5 Conclusions

The *reconstruction of cancer progression models* is a pressing problem, as it promises to highlight important clues about the evolutionary dynamics of tumors and to help in better targeting therapy to the tumor (see e.g. Olde Loohuis et al., 2014). In the absence of large longitudinal datasets, progression extraction algorithms rely primarily on *cross-sectional* input data, thus complicating the statistical inference problem.

In this paper we presented CAPRI, a new algorithm (and part of the TRONCO package) that attacks the progression model



**Fig. 5.** Atypical chronic myeloid leukemia. (left) Mutational profiles of  $n=64$  ACML patients—exome sequencing in Piazza et al. (2013)—with alterations in  $|G|=9$  genes with either mutation frequency  $> 5\%$  or belonging to an hypothesis input to CAPRI (Supplementary information Section 4). Mutation types are classified as *nonsense point*, *missense point* and *insertion/deletions*, yielding  $m=16$  input events. Purple annotations report the frequency of mutations per sample. (right) Progression model inferred by CAPRI in supervised mode. Node size is proportional to the marginal probability of each event, edge thickness to the confidence estimated with 1000 non-parametric bootstrap iterations (numbers shown leftmost of every edge). The  $P$ -value of the hypergeometric test is displayed too. Hard exclusivity patterns input to CAPRI are indicated as red squares. Events without inward/outward edges are not shown

reconstruction problem by inferring *selectivity relationships* among ‘genetic events’ and organizing them in a graphical model. The reconstruction algorithm draws its power from a combination of a scoring function (using Suppes’ conditions) and subsequent filtering and refining procedures, maximum-likelihood estimates and bootstrap iterations. We have shown that CAPRI outperforms a wide variety of state-of-the-art algorithms. We note that CAPRI performs especially well in the presence of noise in the data, and with limited sample size. Moreover, we note that, unlike other approaches, CAPRI can reconstruct different types of confluent trajectories unaffected by the irregularities in the data—the only limitation being our ability to hypothesize these patterns in advance. We also note that CAPRI’s overall algorithmic complexity and convergence properties do offer several tradeoffs to the user.

Successful cancer progression extraction is complicated by tumor heterogeneity: many tumor types have molecular subtypes following different progression patterns. For this reason, it can be advantageous to cluster patient samples by their genetic subtype prior to applying CAPRI. Several tools have been developed that address this clustering problem (e.g. Network-based stratification Hofree *et al.*, 2013 or COMET from Leiserson *et al.*, 2015). A related problem is the classification of mutations into functional categories. In this paper, we have used genes with deleterious mutations as driving events. However, depending on other criteria, such as the level of homogeneity of the sample, the states of the progression can represent any set of discrete states at varying levels of abstraction. Examples include high-level hallmarks of cancer proposed by Hanahan and Weinberg (2000, 2011), a set of affected pathways, a selection of driving genes, or a set of specific genomic aberrations such as genetic mutations at a more mechanistic level.

We are currently using CAPRI to conduct a number of studies on publicly available datasets (mostly from TCGA, NCI and the NHGRI, 2005) in collaboration with colleagues from various institutions. In this work we have shown the results of the reconstruction on the ACML dataset published by Piazza *et al.* (2013), and in Supplementary information Section 4 we include a further example application on ovarian cancer (Knutsen *et al.*, 2005), as well as a comparative study against the competing techniques. Furthermore, we are currently extending our pipeline to include pre-processing functionalities, such as patient clustering and categorization of mutations/genes into pathways (using databases such as the KEGG database (see Kanehisa and Goto, 2000) and functionalities from tools like Network-based clustering, due to Hofree *et al.* (2013).

Encouraged by CAPRI’s ability to infer interesting relationships in a complex disease such as aCML, we expect that in the future CAPRI will help uncover relationships to aid our understanding of cancer and eventually improve targeted therapy design.

## Acknowledgements

We also thank Francesca Ciccarelli, King’s College London, UK, and others for suggesting the ‘selectivity advantage’ terminology. We would also like to thank all the participants of the Workshop and School on *Cancer, Systems and Complexity* held on Lake Como, Italy for many fruitful discussions there (csac.lakecomoschool.org). Finally, we are also indebted to Rocco Piazza, Università degli Studi di Milano Bicocca, Italy, for all the data, insights and patience in explaining to us the biology of aCML.

## Funding

This research was funded by the NSF grants CCF-0836649 and CCF-0926166 and by Regione Lombardia (Italy) under the research projects RetroNet through the ASTIL Program [12-4-5148000-40]; U.A. 053 and

Network Enabled Drug Design project [ID14546A Rif SAL-7] Fondo Accordi Istituzionali 2009.

*Conflict of Interest:* none declared.

## References

- Abdel-Wahab, O. *et al.* (2012) Asx1 mutations promote myeloid transformation through loss of prc2-mediated gene repression. *Cancer Cell*, **22**, 180–193.
- Antonioti, M. *et al.* (2014) The TRONCO package for translational oncology. Available at standard R repositories.
- Attolini, C.S.-O. *et al.* (2010) A mathematical framework to determine the temporal sequence of somatic genetic events in cancer. *Proc. Natl. Acad. Sci.*, **107**, 17604–17609.
- Beerenwinkel, N. *et al.* (2005) Learning multiple evolutionary pathways from cross-sectional data. *J. Comput. Biol.*, **12**, 584–598.
- Beerenwinkel, N. *et al.* (2007) Conjunctive bayesian networks. *Bernoulli*, **13**, 893–909.
- Beerenwinkel, N. *et al.* (2014) Cancer evolution: mathematical models and computational inference. *Syst. Biol.*, **64**, e1–e25.
- Carvalho, A.M. (2009) Scoring functions for learning Bayesian networks. *Inesc-id Tec. Rep.*
- Cheng, Y.-K. *et al.* (2012) A mathematical methodology for determining the temporal order of pathway alterations arising during gliomagenesis. *PLoS Comput. Biol.*, **8**, e1002337.
- Desper, R. *et al.* (1999) Inferring tree models for oncogenesis from comparative genome hybridization data. *J. Comput. Biol.*, **6**, 37–51.
- Desper, R. *et al.* (2000) Distance-based reconstruction of tree models for oncogenesis. *J. Comput. Biol.*, **7**, 789–803.
- Efron, B. (1982) *The Jackknife, the Bootstrap and Other Resampling Plans*, Volume 38 of CBMS-NSF Regional Conference Series in Applied Mathematics. SIAM, Philadelphia, Pennsylvania, USA.
- Efron, B. (2010) *Large-scale inference: empirical Bayes methods for estimation, testing, and prediction*, Volume 1. Cambridge University Press, Cambridge, UK.
- Figueroa, M.E. *et al.* (2010) Leukemic idh1 and idh2 mutations result in a hypermethylation phenotype, disrupt tet2 function, and impair hematopoietic differentiation. *Cancer Cell*, **18**, 553–567.
- Gelsi-Boyer, V. *et al.* (2009) Mutations of polycomb-associated gene asx1 in myelodysplastic syndromes and chronic myelomonocytic leukaemia. *Br. J. Haematol.*, **145**, 788–800.
- Gerstung, M. *et al.* (2009) Quantifying cancer progression with conjunctive bayesian networks. *Bioinformatics*, **25**, 2809–2815.
- Gupta, A. and Bar-Joseph, Z. (2008) Extracting dynamics from static cancer expression data. *IEEE/ACM Trans. Comput. Biol. Bioinf.*, **5**, 172–182.
- Hanahan, D. and Weinberg, R.A. (2000) The hallmarks of cancer. *Cell*, **100**, 57–70.
- Hanahan, D. and Weinberg, R.A. (2011) Hallmarks of cancer: the next generation. *Cell*, **144**, 646–674.
- Heckerman, D. *et al.* (1995) Learning bayesian networks: the combination of knowledge and statistical data. *Mach. Learn.*, **20**, 197–243.
- Hitchcock, C. (2012) Probabilistic causation. In: Zalta, E.N. (ed.) *The Stanford Encyclopedia of Philosophy*. Stanford University, winter 2012 edition.
- Hjelm, M. *et al.* (2006) New probabilistic network models and algorithms for oncogenesis. *J. Comput. Biol.*, **13**, 853–865.
- Hofree, M. *et al.* (2013) Network-based stratification of tumor mutations. *Nat. Methods*, **10**, 1108–1115.
- Huang, S. *et al.* (2009) Cancer attractors: a systems view of tumors from a gene network dynamics and developmental perspective. *Semin. Cell Dev. Biol.*, **20**, 869–76.
- Inoue, D. *et al.* (2014) Setbp1 mutations drive leukemic transformation in asx1-mutated mds. *Leukemia*, **29**, 847–857.
- Kanehisa, M. and Goto, S. (2000) Kegg: kyoto encyclopedia of genes and genomes. *Nucleic Acids Res.*, **28**, 27–30.
- Kleinberg, S. (2012) *Causality, Probability, and Time*. Cambridge University Press, Cambridge, UK.

- Knutsen, T. et al. (2005) The interactive online sky/m-fish & cgh database and the entrez cancer chromosomes search database: Linkage of chromosomal aberrations with the genome sequence. *Genes Chromosomes Cancer*, **44**, 52–64.
- Koller, D. and Friedman, N. (2009) *Probabilistic Graphical Models: Principles and Techniques—Adaptive Computation and Machine Learning*. The MIT Press, Cambridge, Massachusetts, USA.
- Leiserson, M. et al. (2015) Comet: A statistical approach to identify combinations of mutually exclusive alterations in cancer. In: Proceedings of the 19th Annual Research in Computational Biology Conference (RECOMB).
- Lin, C.-C. et al. (2014) Sf3b1 mutations in patients with myelodysplastic syndromes: The mutation is stable during disease evolution. *Am. J. Hematol.*, **89**, E109–E115.
- Magwene, P.M. et al. (2003) Reconstructing the temporal ordering of biological samples using microarray data. *Bioinformatics*, **19**, 842–850.
- Meggendorfer, M. et al. (2013) Setbp1 mutations occur in 9% of mds/mpn and in 4% of mpn cases and are strongly associated with atypical cml, monosomy 7, isochromosome i (17)(q10), asxl1 and cbl mutations. *Leukemia*, **27**, 1852–1860.
- Merlo, L.M. et al. (2006) Cancer as an evolutionary and ecological process. *Nat. Rev. Cancer*, **6**, 924–935.
- Misra, N. et al. (2014) Inferring the paths of somatic evolution in cancer. *Bioinformatics*, **30**, 2456–2463.
- Muto, T. et al. (2013) Concurrent loss of ezh2 and tet2 cooperates in the pathogenesis of myelodysplastic disorders. *J. Exp. Med.*, **210**, 2627–2639.
- NCI and the NHGRI. (2005) The Cancer Genome Atlas.
- Olde Loohuis, L. et al. (2014) Cancer hybrid automata: model, beliefs & therapy. *Inf. Comput.*, **236**, 68–86.
- Olde Loohuis, L. et al. (2014) Inferring tree causal models of cancer progression with probability raising. *PloS one*, **9**, e115570.
- Pearl, J. (1988) *Probabilistic reasoning in intelligent systems: networks of plausible inference*. Morgan Kaufmann, San Francisco, California, USA.
- Pearl, J. (2000) *Causality: models, reasoning and inference*, volume **29**. Cambridge Univ Press, Cambridge, UK.
- Piazza, R. et al. (2013) Recurrent setbp1 mutations in atypical chronic myeloid leukemia. *Nat. Genet.*, **45**, 18–24.
- Schwarz, G. (1978) Estimating the dimension of a model. *Ann. Stat.*, **6**, 461–464.
- Spirtes, P. et al. (2000) *Causation, prediction, and search*, Volume **81**. MIT press, Cambridge, Massachusetts, USA.
- Suppes, P. (1970) *A Probabilistic Theory of Causality*. North-Holland Publishing Company, Amsterdam, Holland.
- Szabo, A. and Boucher, K. (2002) Estimating an oncogenetic tree when false negatives and positives are present. *Math. Biosci.*, **176**, 219–236.
- Tamborero, D. et al. (2013) Comprehensive identification of mutational cancer driver genes across 12 tumor types. *Sci. Rep.*, **3**, 1–9.
- Tsamardinos, I. et al. (2003) Algorithms for large scale markov blanket discovery. In: FLAIRS Conference, Volume **2003**, p. 376–381.
- Vogelstein, B. et al. (1988) Genetic alterations during colorectal-tumor development. *New Engl. J. Med.*, **319**, 525–532.
- Vogelstein, B. et al. (2013) Cancer genome landscapes. *Science*, **339**, 1546–1558.