

PRODIGY: personalized prioritization of driver genes

Gal Dinstag and Ron Shamir*

School of Computer Science, Tel Aviv University

{rshamir,galdinstag}@tau.ac.il

* corresponding author

Date: 30.10.2018

Abstract

Background: Evolution of cancer is driven by few somatic mutations that disrupt cellular processes, causing abnormal proliferation and tumor development, while most somatic mutations have no impact on progression. Distinguishing those mutated genes that drive tumorigenesis in a patient is a primary goal in cancer therapy: Knowledge of these genes and the pathways on which they operate can illuminate disease mechanisms and indicate potential therapies and drug targets. Current research focuses mainly on cohort-level driver gene identification, but patient-specific driver gene identification remains a challenge.

Methods: We developed a new algorithm for patient-specific ranking of driver genes. The algorithm, called PRODIGY, analyzes the expression and mutation profiles of the patient along with data on known pathways and protein-protein interactions. Prodigy quantifies the impact of each mutated gene on every deregulated pathway using the prize collecting Steiner tree model. Mutated genes are ranked by their aggregated impact on all deregulated pathways.

Results: In testing on five TCGA cancer cohorts spanning >2500 patients and comparison to validated driver genes, Prodigy outperformed extant methods and ranking based on network centrality measures. Our results emphasize the pleiotropic effect of driver genes and show that Prodigy is capable of identifying even very rare drivers. Hence, Prodigy can assist oncologists in decisions regarding personalized treatment.

Availability: The Prodigy software is available from the authors upon request.

Introduction

Cancer is an evolutionary process in which normal cells accumulate genomic and epigenomic alterations of various kinds including Single Nucleotide Variations (SNVs) and chromosomal aberrations. Some of these alterations confer growth and positive selection advantage to the mutated cells, giving rise to intensive proliferation and tumors¹. The alterations can be inherited through germ line mutations as in the case of *BRCA1* and *BRCA2* in breast cancer² or occur somatically¹. While somatic mutations also occur in normal cells, they are neutral or cause apoptosis but do not lead to transformation into a cancer cell.

Driver mutations: Mutational events that grant such advantages to the cell and "drive" it into tumorigenesis are called *driver mutations* (or driver events) and the genes in which these mutations take place are called *driver genes*. In contrast, *passenger mutations* are acquired extensively during cancer progression simply because cancer cells over-proliferate in orders of magnitude as compared to normal cells, and random mutations mainly occur during cell division. These mutations do not confer any growth advantage³.

The overall number of observed mutations varies among tumor tissues. Kim & Kim⁴ analyzed dozens of cancer patient cohorts from TCGA⁵ and found that the average number of somatic mutations can reach up to thousands per tumor in some cancer subtypes. There is a very extensive debate regarding the number of driver mutations among the observed mutations in each tumor^{1,3,6} but the consensus is that this number is very low. Obviously, there are many factors that contribute to the variation in the number of drivers, including the progression stage of the tumor⁷, its tissue of origin⁸, environmental properties such as smoking⁹ and other factors like age¹⁰. Tomasetti et al.¹¹ showed that as little as three driver mutations suffice to develop lung and colorectal cancer. Nordling¹² and Armitage¹³ suggested six or seven as the typical number of drivers.

It is therefore a challenge to distinguish driver from passenger mutations. The need to do so has high priority in cancer research - and in personalized cancer medicine in particular - for several reasons: 1) knowledge of the drivers and the mechanisms by which they operate can suggest potential treatments and drug targets. 2) Basing cancer treatment on molecular signatures rather than on the disease organ offers the opportunity to treat individuals with regimens not yet considered for their specific type of cancer. For example, many "basket" clinical trials, in which a specific drug is given to patients with diverse cancer types based on specific biomarkers, show that the same drug can have high efficiency across different types if the right mutation is detected¹⁴.

Driver gene identification in large cohorts: Computational research regarding driver genes first focused on distinguishing driver mutations from passengers in a cohort of patients (usually of the same tissue of origin): MuSiC¹⁵ uses the statistical significance of higher than expected rate of mutations, along with pathway mutation rate and correlation with clinical features to detect drivers. MutSigCV¹⁶ estimates the background mutation rate of each gene and identifies mutations that significantly deviate from that rate. MEMo¹⁷ tries to find small subnetworks of genes that belong to the same pathway and exhibit internal mutual exclusivity patterns. HotNet2¹⁸ incorporates knowledge from protein-protein interaction (PPI) networks to find small subnetworks of frequently mutated genes using heat-diffusion process. TieDie¹⁹ also incorporates PPIs and mRNA expression data to find overlapping subnetworks that possess high degree of mutation and expression values using heat-diffusion. DriverNet²⁰ tries to find a parsimonious set of mutated genes that is linked to genes that experience deregulation of mRNA expression in a given PPI network. Paradigm-Shift²¹ utilizes SNV, copy number variation (CNV), expression and known pathways to infer gain or loss-of-function of mutated genes in single patients. Many more methods for driver gene detection in cohorts are covered by Chang et al.²² and Tokahim et al.²³.

The methods above focus on general driver gene detection, but do not aim to offer personalized means of diagnosis or treatment: individual patients may have different compositions of mutated driver genes (**Supp. Fig. 1**). In addition, these methods rely on statistical power obtained by large cohorts and by

doing so, they inevitably underestimate the importance of rare drivers that occur in only a handful of patients (also known as the "long tail phenomenon"²⁴) and are important only for them. Here we focused on patient specific driver prioritization.

Although many driver mutations were experimentally validated³⁶, personalized driver prioritization is needed for several reasons: 1) Some patients carry mutations in dozens of known drivers (**SFig 1**), and it is essential to understand which are the true drivers for the patient. 2) Some patients do not possess mutations in any known driver (**SFig 1**), so one has to find putative ones de novo. 3) Even if a patient has only few mutations in known drivers, and assuming they are all active, we still need to internally rank them, since the number of therapies that can be given to an individual at the same time is very low due to toxicity and adverse events^{25,26}.

Personalized driver gene profiles: To address the need for personalized driver gene identification and prioritization, one must develop methods that can operate on the data of a single patient. Several attempts have been made in this direction: DawnRank²⁷ uses a variant of Google's PageRank to rank an individual's mutated genes profile according to its effect on expression deregulation of downstream genes in a large directed PPI network. It ranks the genes by quantifying the impact of each of them on the differentially expressed genes (DEGs) using a diffusion process. SCS²⁸ tries to find a parsimonious set of mutated genes that are linked to downstream DEGs in a large directed PPI network. These methods rank putative driver genes for a patient. In contrast, Hitn'DRIVE²⁹ outputs a set of candidate driver genes without internal ranking. It tries to find a parsimonious set of mutations with short expected path lengths to a set of DEGs. The lack of ranking is a drawback from a treatment perspective, especially when the number of predicted genes is large.

This study: Here we develop a new algorithm for ranking of driver genes of an individual. The algorithm, called PRODIGY (Personalized Ranking Of DrIver Genes analYsis) scores mutations by their influence on deregulation of multiple known pathways. Unlike the methods described above, Prodigy collects multiple signals from many local views of the same tumor rather than one global view. These local views are based on curated pathways and each one reflects a different aspect of the deregulation state of the tumor. Thus, the extent to which a given mutation explains multiple pathway deregulations serves as a proxy to the likelihood that this mutation is indeed one of the drivers. Our algorithm assumes that driver mutations influence the deregulation of other genes in affected pathways. In particular the true drivers will have good connectivity to these pathways, and our method is designed to score such connections correctly using a variant of the prize collecting Steiner tree problem. By aggregating many local views for all mutations of an individual, a global picture can be made and the personalized landscape of drivers can be assembled and ranked.

In testing on five TCGA cancer cohorts spanning >2500 patients and comparison to validated driver genes, PRODIGY outperformed extant methods and ranking based on network centrality measures. Our results emphasize the pleiotropic effect of driver genes and show that PRODIGY is capable of identifying even very rare drivers. Hence, PRODIGY can assist oncologists in decisions regarding personalized treatment.

Caveats: Note that while we occasionally talk about driver mutations, all our analysis is done on the gene level and - as in SCS and DawnRank - different mutations in the same gene are not distinguished. Since the number of mutations per mutated gene in a patient is usually 1 (**Supp. Table 1**) this distinction is less important for personalized ranking than for cohort-level analyses. Also, as we shall see, often we identify and rank ten genes or more per patient, so the notion of drivers in this study is somewhat more lenient than is common in the literature. However, our results suggest that a larger number of predicted drivers actually contribute to the performance.

Methods

Given the set of mutated genes and the expression profile of an individual, we wish to rank the mutated genes in that individual. Our assumption is that the influence of driver genes is disseminated

along pathways and is manifested by DEGs. By aggregating evidence from multiple pathways for a mutated gene, we score the extent to which it explains the deregulation of the pathways. This score serves as a proxy to the likelihood that the gene is a driver in the patient. Mathematically, we score the influence of a mutation on a deregulated pathway using the undirected prize collecting Steiner tree (PCST) model.

The PCST model: In this problem (**Figure 1A**) the goal is to find in a weighted graph a subtree maximizing the sum of the weights of the nodes minus the cost of edges in it. The input is an undirected graph $G = (V, E, W, P)$. $W: E \rightarrow R_+$ is a positive weight function on the edges and $P: V \rightarrow R$ is a weight function on the nodes. In our context, edge weights are penalties reflecting interaction reliability, positive node weights are prizes given to DEGs as they reflect the pathway deregulation that we want to capture, while other nodes that can serve as intermediate nodes in the tree (*Steiner nodes*) are assigned non-positive values serving as penalties. Given a node $g \in V$, the objective is to find a subtree T of G that contains g and maximizes:

$$\text{Score}(T) = \sum_{v \in V_T} P(v) - \sum_{(u,v) \in E_T} W(u, v)$$

In other words, the score of T is the total profit of pre-defined prizes minus the penalties of using intermediate edges and nodes. This model was shown to be suitable in different biological problems and in particular in scenarios where a mechanistic view is desirable^{30,31}.

Data and reference network: Prodigy uses two types of genomic data for each patient: the list of mutated genes, i.e. all genes with SNVs or small insertions/deletions in coding regions, and the profile of mRNA expression. mRNA expression profiles from healthy tissue samples are also utilized for differential expression analysis. Prodigy also uses two types of undirected interaction networks: 1) a global PPI network taken from STRING v10.5³². Here we used only physical interactions that were validated experimentally and interactions from other curated databases with confidence score > 0.7 , so that only highly reliable interactions were included. The resulting network had 11,302 nodes and 273,210 edges. 2) A collection of pathways. Here we used either Reactome³³, NCI PID³⁴ or KEGG³⁵. Information about the pathway databases is given in **STable 2**.

The Prodigy algorithm

A schematic view of the algorithm is given in **Figure 1**. The algorithm works as follows:

Pre-processing Given a patient's mRNA expression profile (as read counts), differential expression analysis was done using DeSEQ2³⁶ by comparing the profile to a background expression distributions from healthy samples of the same tissue of origin. All genes with > 2 log2-fold-change that are statistically significant (FDR = 0.05) were identified as DEGs.

The gene set of each pathway is tested for enrichment in DEGs using the hyper-geometric score, and pathways that are significantly enriched (FDR = 0.05) are called *deregulated*.

Driver - pathway scores We use a global interaction network $G = (V, E, W)$ where W is the edge confidence score. For a deregulated pathway p we also have its network $G_p = (V_p, E_p)$. Both networks are undirected. The influence score of the mutated gene g on pathway p is calculated as follows:

1. We construct a new network $G_{p,g} = (V_{p,g}, E_{p,g}, W_{p,g}, P_{p,g})$ that is derived from G, G_p and g , as follows: The nodes of the network are those of the deregulated pathway, g , and $N(V_p \cup g)$ - their distance 1 neighbors in G :

$$V_{p,g} = V_p \cup g \cup N(V_p \cup g)$$

Its edges are those of the deregulated pathway plus all edges of the global network with both ends in $V_{p,g}$:

$$E_{p,g} = E_p \cup \{(u, v) | u, v \in V_{p,g} \text{ and } (u, v) \in E\},$$

The cost of the edges from p is 0.1. For the other edges, which originate from the global network G , their cost depends on their confidence score in that network, with edges of higher confidence costing less.

$$W_{p,g}(u, v) = \begin{cases} 0.1, & (u, v) \in E_p \\ 1 - W(u, v), & \text{otherwise} \end{cases}$$

Edges from the pathway are assigned a constant penalty of 0.1 since pathway databases do not provide confidence scores for the interactions, but those pathways are highly curated. In contrast, the confidence scores on the edges from the global network are given an upper bound of 0.8 so that their cost in $G_{p,g}$ is at least 0.2. The rationale is that we want to steer the algorithm to prefer the original pathway edges, while allowing some alterations.

Finally every DEG that belongs to the pathway has a positive (prize) score depending on its fold change (FC), and every other node v has a negative (penalty) score depending on its degree in $G_{p,g}$ as follows:

$$P_{p,g}(v) = \begin{cases} \log(FC(v)), & v \in DEG \cap V_p \\ -degree(v)^\alpha, & \text{otherwise} \end{cases}$$

Note that DEGs in $V_{p,g} \setminus V_p$ have negative values. The PCST problem aims to collect as much of the prize nodes value while paying least penalty for intermediate edges and nodes. Intermediate nodes that have high degree ("hubs") open more connection options and are thus penalized higher depending on their degree. The α parameter controls that penalty.

2. Having constructed $G_{p,g}$ we now seek a tree $T_{p,g}$ that contains g of optimal score. If $\text{Score}(T_{p,g}) \leq 0$ (i.e., no tree with positive score is found), an empty tree with score 0 is output instead.
3. To account for variability in pathway size and the number of DEGs in the pathway, the *influence score* of mutated gene g on pathway p is defined as the fraction of attained score out of the upper bound of all positive prizes in the pathway:

$$\text{Infl}(p, g) = \frac{\text{Score}(T_{p,g})}{\sum_{v \in V_{p,g}} \max\{P_{p,g}(v), 0\}}$$

The overall *influence score* of g is $\text{infl}(g) = \sum_{p \in DP} \text{infl}(p, g)$ where DP is the set of deregulated pathways of the patient.

Pathway filtering We compute driver-pathway influence scores for all mutated genes and all deregulated pathways. For the final score we exclude pathways for which more than half of the genes had a positive score. These are mainly very large pathways that have high connectivity in the global network, and therefore some genes may acquire positive influence scores by chance.

Gene filtering Genes that acquired positive scores in many pathways have greater chance to represent a true effect on the tumor than genes that attained positive scores for only few pathways, possibly due to the topology of the network. In some patients, when plotting the distribution of $\text{Infl}(g)$ scores across all mutated genes g (after filtering pathways), we observed a bimodal distribution (see **SFig. 2**). Typically, one distribution contains genes with high scores collected from many pathways and the other contains genes with low scores collected from a few pathways. We modeled this distribution as a mixture of two Gaussians and computed its maximum likelihood parameters using EM³⁷. We then excluded all genes that had higher posterior probability to come from the distribution with the lower mean (**SFig. 2**). In case a bimodal distribution was not observed, we did not filter any gene.

Final ranking After the filtering steps, genes are ranked according to their overall influence scores.

Comparison to other methods:

We compared Prodigy to DawnRank²² and SCS²³. Since both DawnRank and SCS use directed graphs, the global PPI network used to test them was taken from the original publication. This network contained 11,648 nodes and 211,794 directed edges. To ensure that results are not derived primarily from the topology of the network, we also generated personalized rankings using three node centrality measures: node degree, closeness and betweenness (see **Supp. Methods** for definitions). To produce rankings based on each measure, we calculated it on each of the networks $G_{p,g}$ and summed the results over all the networks for a final ranking.

Validation:

In order to validate rankings, we used a curated list of driver genes from the Cancer Gene Census (CGC) as gold standard. CGC is part of COSMIC³⁸, the largest database of somatic mutations in cancer. CGC contains mutations of different forms (gene amplifications, SNVs, translocations etc.) that were experimentally validated as driver mutations for different cancer types. Since we only used information about SNVs and short indels of each patient, we used as ground truth only genes that were classified by CGC as containing a driver SNV or indel ($n = 248$ out of 567). In this validation, we assumed that if a gold-standard gene was mutated in a patient, it is a true driver gene in the patient's tumor. We measured the quality of each method by means of precision, recall and F1 with respect to the gold standard (see **Supp. Methods**)

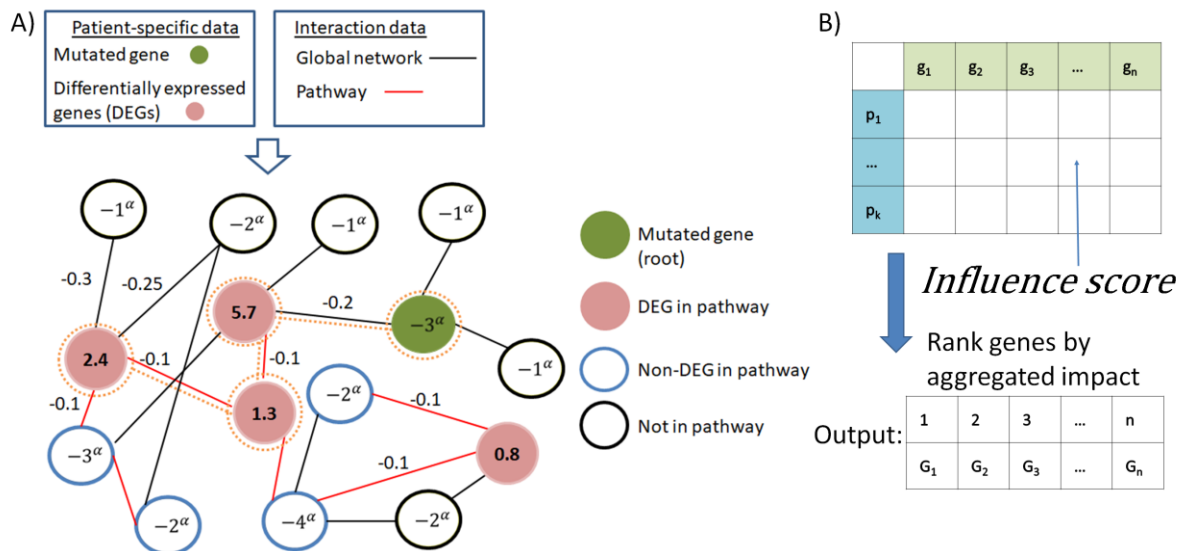


Figure 1: Outline of Prodigy's approach. A) Scoring the influence of the mutated gene g on the pathway p : The pathway and genes at distance 1 from it or from the mutated gene g in the global network, along with the global edges among them, constitute the network $G_{p,g}$ for analysis (see **Methods**). This is the network shown here. Node prizes (positive values) reflect the extent of differential expression of DEGs in p , and node penalties reflect other node's degrees (calibrated by the exponent α). Edge penalties reflect interaction confidence. The goal is to find a maximum weight subtree in the network rooted at the mutated gene g . Its weight is the score of the PCST solution. In this example the subtree marked by orange dotted lines is the PCST solution, of score $9 \cdot 3^\alpha$. The *influence score* of the pair (p, g) is the score of the PCST solution, divided by the sum of the values of DEGs that belong to p (10.2 here). B) After calculating the influence score for all pairs (p, g) , we filter out some pathways and genes from the scoring matrix (see **Methods**). The final output is a ranking of the remaining genes by their aggregated score on the remaining pathways.

Driver-Pathway linkage:

Prodigy can quantify driver-pathway associations, allowing us to explore novel interactions and even cancer subtype-specific ones. Our hypothesis was that if driver gene g often deregulates pathway p then they will be observed together more frequently in patients of the cohort, and the deregulation state of p will be higher when g is acting as a driver. To test this conjecture, we focused on the ten top

ranked genes for each individual and looked for driver-pathway pairs where the number of patients for whom the gene was ranked high and the pathway was deregulated was unexpectedly high according to the hyper-geometric distribution. For each pair, we then tested if p was more deregulated when g was classified as driver using t-test (see **Supp. Methods** and **SFig. 4** for more details).

Results

Driver gene ranking: We tested six ranking methods on 2569 samples from five cohorts of cancer patients from TCGA: COAD, LUAD, BRCA, HNSC and BLCA³⁹⁻⁴³ (212, 487, 969, 502 and 399 samples, respectively). We used a training set comprised of 10% of the samples from each cohort to derive the optimal node degree weighting factor α in terms of F1, and used the chosen parameter to calculate personalized rankings for the remaining 90%. Prodigy's results were consistent across different α values (**SFig. 4**) with significant decline in performance for values > 0.2 . $\alpha = 0.05$ was chosen for all cohorts.

Figure 2A shows the average precision, recall and F1 for Prodigy, DawnRank and for the three centrality measures using the Reactome pathways (see **Methods**). The results are reported as average values for the entire cohort as a function of the top N ranked genes. If an individual had less than N ranked genes, the last value for this patient was duplicated so that all quality measure vectors for all patients are of length N. Since SCS reported empty rankings for 720 samples (28%), it is not shown in **Figure 2A**. Performance of all methods on the set of patients for whom SCS produced results (the "SCS sub-cohort") is shown in **SFig. 6**, and performance for different cancer types is shown in **SFig.7-9**. Results for the KEGG and NCI pathway databases for the entire cohort were similar (**SFig. 5**).

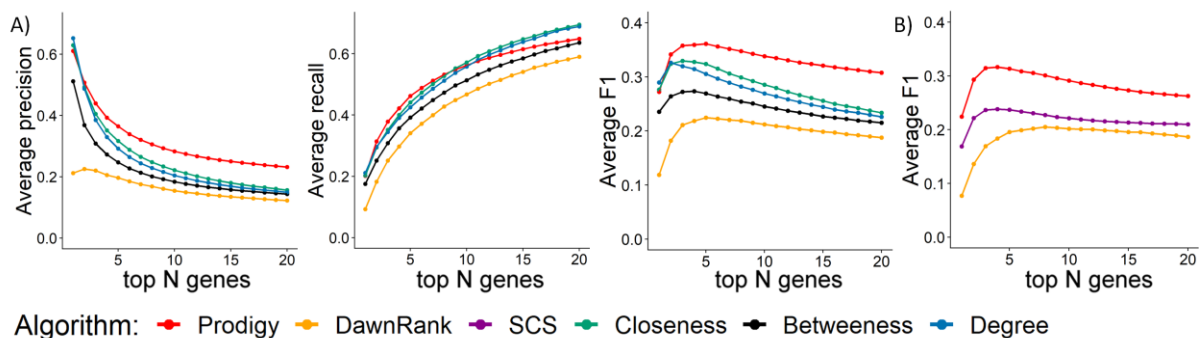


Figure 2: A. Average precision, recall and F1 across all patients ($n=2340$) as a function of the number of top ranked genes in the personalized profiles. Prodigy's results were derived using the STRING global PPI network (see **Methods**) and Reactome pathways B. Average F1 using the global network from the SCS and DawnRank papers, on those patients for whom SCS proposed drivers ($n=1804$).

Overall, Prodigy outperformed SCS and DawnRank in terms of F1, precision and recall. On the NCI pathways and for high values of N on KEGG pathways, SCS was better for the SCS sub-cohort. To ensure that the improvement in results does not stem from the difference in the underlying networks, we also tested Prodigy on the same network used by DawnRank and SCS with two adjustments: (1) Since Prodigy works on undirected graphs, we ignored edge directions. (2) Since this network is unweighted, we gave weight=0.2 to all edges (and 0.1 to pathway edges as before, see **Methods**). The results (**Figure 2B** and **SFig. 10**) clearly show that Prodigy outperforms DawnRank and SCS even on their network.

Remarkably, the centrality measures produced very good predictions, consistently better than DawnRank and SCS – but worse than Prodigy. These measures had better recall than Prodigy, probably due to the fact that no filtering was done on the centrality measures while Prodigy excluded genes not likely to be drivers for an individual. The fact that driver genes are associated with high network connectivity was previously discussed^{29,44,45} and we observed it as well: in our global

network derived from STRING, known drivers included in the CGC tended to have high degree and betweenness (SFig. 8). Our results emphasize the need to account for "hubness" property in methods for driver gene ranking. Prodigy accounts for this factor by penalizing Steiner nodes according to their degree. Taken together the results clearly demonstrate that Prodigy outperforms mere topology measures in capturing true driver genes.

Discovering rare drivers: One of the advantages of Prodigy is its ability to identify rare drivers, even when the gene is mutated in few patients. To demonstrate this ability we looked for mutated genes that had frequency < 2% in the cohort and were ranked in the top 10 drivers of individuals. The results are summarized in Figure 3. In some cohorts, most of the mutated genes were in fact rare (< 2%, STable 3), which is of course reflected in our results. On the other hand, Prodigy prioritized rare mutations to lesser extent than their frequency in the population (STable 3).

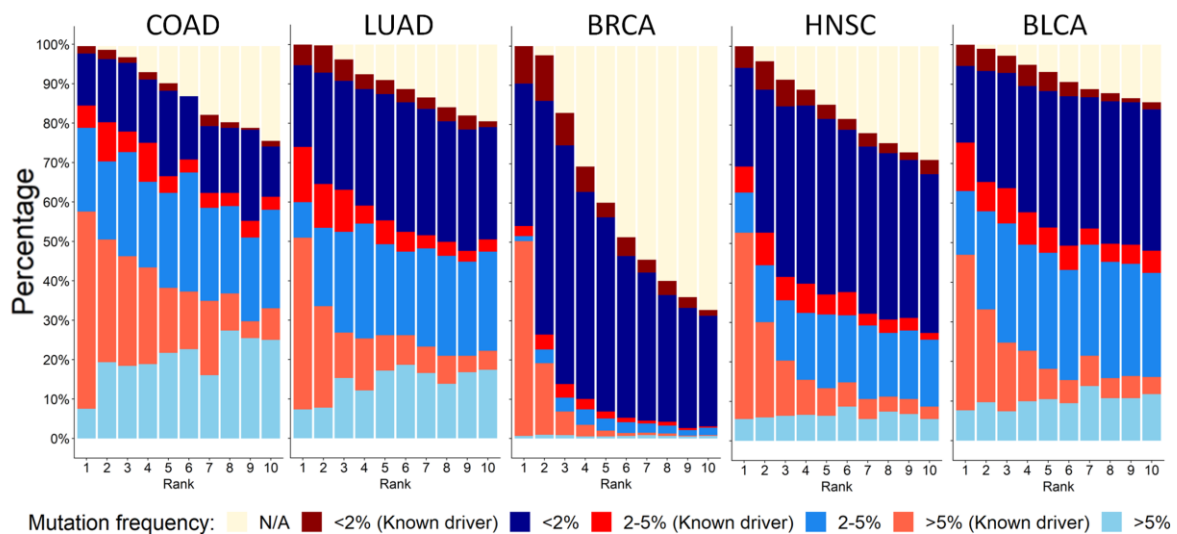


Figure 3: Prodigy discovers rare drivers. For each cancer type and for each individual we analyzed the top 10 genes according to the ranking. For $k = 1, \dots, 10$ the plot shows the fraction of patients for whom the gene ranked k -th belongs to the respective frequency bin (as denoted by its color). N/A: patients for whom Prodigy ranked less than k mutations.

Prodigy was able to detect known rare drivers. For example, for the colon cancer patient *TCGA-AD-6899*, Prodigy ranked highly the gene *SRC*, a known driver in colon cancer. Remarkably, this patient was the only one (out of 212) who harbored a mutation in that gene. In HNSC, *FES* mutation was observed in five patients out of 502 (1%), and was highly ranked in two of them. *MTOR* was mutated in nine patients (1.7%) and highly ranked in five. *TSC2* was mutated in six patients (1.1%) and highly ranked in two. All of these genes are known HNSC-specific drivers according to CGC. In LUAD, *HIF1A* was mutated in two patients and was highly ranked in both. *RAD21* was mutated in eight (1.6%) and highly ranked in one. *ARAF* was mutated in five (1%) and highly ranked in one. *EED* was mutated in six and highly ranked in one. These are all known drivers of LUAD. The results show that Prodigy is capable of identifying even very rare drivers from the CGC. Taken together, we demonstrated Prodigy's ability to detect both rare and frequent drivers.

Driver gene-pathway linkage: We identified 1299 significant driver-pathway interactions (see [Suppl. File 1](#)). They include some very well-known interactions between *TP53* and sub pathways of the cell cycle in all cohorts except COAD and *TP53*-DNA repair pathways in the BLCA cohort. Moreover, the gene *A2M* was associated with "G alpha (i) signaling events" in the COAD, BRCA and BLCA cohorts. The G alpha (i) signaling pathway belongs to the GPCR family of signaling pathways, which are strongly linked to cancer⁴⁶. This analysis can provide new insights on the mechanism by which the drivers operate and can offer new targets for further research.

Multi-pathway effect: One of the main assumptions underlying Prodigy is that driver genes affect cellular process pathways, and therefore summarized scores from multiple pathways will improve our ability to identify them. This is in contrast to previous methods that took a global approach to driver gene prioritization based on a single unified picture of the state of the tumor^{22,23}. In order to test whether multiple sources indeed contribute to the accuracy of prediction, we explored the performance as a function of the number k of allowed pathways per mutated gene. For $k = 1, \dots, 50$, we used the top k scoring pathways of each gene for ranking and examined the average area under the precision-recall curve (AUPR) for each cohort (see **Supp. Methods**). **Figure 4A** shows that for all cohorts, AUPR improved with incorporating more pathways and plateaued at 15-30 pathways.

Since different pathways may partially overlap, we tested the extent of this overlap and its effect on performance. We computed the distribution of Jaccard Index between pairs in the top 20 scoring pathways of each gene (i.e., the number of genes that belong to both pathways divided by the number of genes in their union, **Supp. Methods**). The results show substantial overlap among the pathways that contribute to the rankings (**Figure 4B**). However, when we filtered out such overlapping pathways, assuming they contain the same information and thus unnecessary for accurate prediction, performance only moderately degraded (**Supp. Methods** and **Figure 4C**). Taken together, we demonstrated the usefulness of using multiple pathways in order to rank driver genes, even when there are overlaps among them.

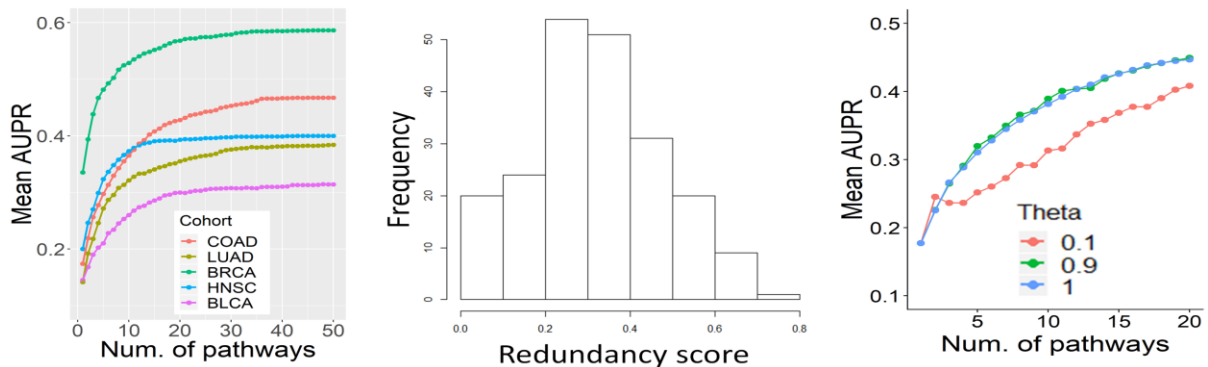


Figure 4: Multi-pathway effect: A) Mean AUPR as a function of the number of top scoring pathways per gene used to derive the results. B) The distribution of redundancy between the top 20 pathways per patient in the COAD cohort ($n = 212$, see **Supp. Methods**). C) Removal of pathway redundancy. The plot shows the AUPR for predicting driver genes in the COAD cohort when filtering out overlapping pathways among the top scoring pathways per gene (**Supp. Methods**). θ is the maximum allowed Jaccard Index between included pathways ($\theta = 1$ implies no filtering).

Actionable and druggable targets: Prodigy's rankings can aid the oncologist in deciding on a patient's therapy, by matching treatment to the predicted driver genes. In order to explore this possibility we used two data sources: (1) DGIdb 3.0⁴⁷, which contains drug targets (or *druggable genes*, i.e., genes with directed pharmacotherapy). Here we used only cancer-specific sources from DGIdb and identified 1375 genes. (2) TARGET⁴⁸, which lists *actionable genes* (i.e., genes for which a genomic-driven therapy exists). The total number of actionable genes was 60. We explored not only the ranked mutated genes themselves but also the pathways that were highly linked (influence score > 0.8) to at least one gene of the top 10 ranked genes of an individual. The rationale is that these pathways are most altered by the driver genes and thus can be targeted in potential treatments. The results (**Figure 5**) indicate that most patients harbor at least one druggable driver (a druggable gene that was prioritized as a driver by Prodigy; mean: 3.32, sd: 2.01) but many do not contain any actionable drivers (mean: 0.82, sd: 0.87). As expected, the number of target genes increased substantially when genes from highly linked pathways were also considered. More importantly, the number of patients without any druggable or actionable gene decreased below 10%. The only exception was the HNSC cohort, where the number of patients without actionable genes remained high (35.8%) even when considering pathways. Hence, Prodigy is able to suggest possible therapeutic

targets personally tailored to the patient's driver genes and uses information about the pathways that are deemed altered by the drivers in order to expand pharmacotherapy options.

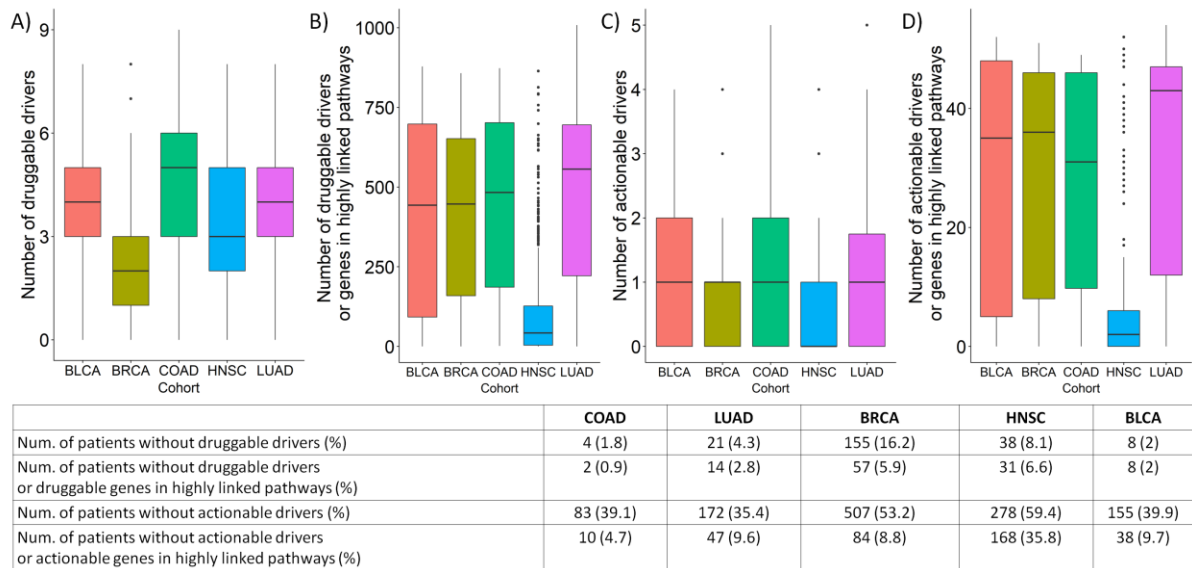


Figure 5: Actionable and druggable genes. The box plots show the distribution of the number of actionable and druggable genes (i.e. genes from TARGET⁴³ and DGIdb⁴²) per patient across the different cohorts. A and C: The distribution of the number of druggable and actionable drivers among the 10 top genes ranked by Prodigy. B and D: The distribution of the number of druggable and actionable genes among predicted drivers and their highly linked pathways. The table describes the number of patients without any druggable/actionable genes in the four categories with respect to the cohort.

Implementation: Prodigy was implemented in R and the software will be publicly available upon publication. The PCST code of⁴⁹ was used. Mean runtime was about 5 minutes per patient on a 65 core, Intel(R) Xeon(R) 2.30GHz, 755GB RAM server.

Discussion:

Personalized diagnosis of cancer patients is a precondition for planning treatment. Deciphering the altered mechanisms and the mutated genes driving them gives a comprehensive picture of the state of the tumor. Although many driver mutations were experimentally validated, there is great potential benefit in identifying which genes act as drivers *in an individual* and prioritizing them: Driver genes are diverse even within a cancer subtype, and they may be rare or not match the disease organ.

Here we provide a novel algorithm called Prodigy for driver gene prioritization in an individual based on the patient's tumor mutation and expression profiles. In testing of over 2500 patients from five cancer subtypes, Prodigy substantially outperformed prior methods for the task. All methods (including ours) use an underlying global interactions network, and we observed that using simple centrality measures of that network to prioritize genes gives better results than the prior methods, but not ours. In all our analyses, as in prior studies, no experimental patient-specific driver information was available, and therefore we used the CGC collection of (global) driver genes as our gold standard.

A unique feature of Prodigy is the fact that it collects information from many distinct pathways in the analysis of a patient. Our results show that driver genes influence many pathways, and that the pathways perspective is more powerful than previous approaches that utilized one global network for the analysis. While Prodigy ranks the genes without setting a cutoff for driver detection, our analysis shows the F1 scores peak around $N=5$. On the other hand, recall rises for $N>5$, so by using prior knowledge about driver genes and observing the actual influence scores of genes that are ranked lower, additional drivers can be pinpointed and used.

Our analysis shows that Prodigy can identify even very rare driver genes, and can reveal linkage between a driver gene and pathways that are preferentially deregulated when the gene acts as a driver. The identified genes typically have multiple drug targets, and thus can suggest treatment decisions.

Acknowledgements:

We thank Nimrod Rappoport for helpful comments on the manuscript. The results published here are based upon data generated by The Cancer Genome Atlas managed by the NCI and NHGRI. Information about TCGA can be found at <http://cancergenome.nih.gov>. Study supported in part by the Israel Science Foundation (grant No. 1339/18), by the DIP German-Israeli Project cooperation, and by Len Blavatnik and the Blavatnik Family foundation. GD was supported by a fellowship from the Edmond J. Safra Center for Bioinformatics at Tel Aviv University.

References:

1. Stratton, M. R., Campbell, P. J. & Futreal, P. A. The cancer genome. *Nature* **458**, 719–724 (2009).
2. King, M.-C., Marks, J. H., Mandell, J. B. & Group, T. Breast and ovarian cancer risks due to inherited mutations in BRCA1 and BRCA2. *Science* (80-.). **302**, 643–646 (2003).
3. Vogelstein, B. *et al.* Cancer genome landscapes. *Science* (80-.). **340**, 1546–1558 (2013).
4. Kim, H. & Kim, Y. M. Pan-cancer analysis of somatic mutations and transcriptomes reveals common functional gene clusters shared by multiple cancer types. *Sci. Rep.* **8**, 1–14 (2018).
5. Weinstein, J. N. *et al.* The Cancer Genome Atlas Pan-Cancer Analysis Project. *Nat. Genet.* **45**, 1113–1120 (2013).
6. Anna C. Schinzel, W. C. H. Oncogenic transformation and experimental models of human cancer. *Front. Biosci.* 71–84 (2008).
7. Vogelstein, B. & Kinzler, K. W. The Path to Cancer- Three Strikes and You're Out. *N. Engl. J. Med.* **373**, 1893–1895 (2015).
8. Kandoth, C. *et al.* Mutational landscape and significance across 12 major cancer types. *Nature* **502**, 333–339 (2013).
9. Govindan, R. *et al.* Genomic Landscape of Non-Small Cell Lung Cancer in Smokers and Never-Smokers. *Cell* **150**, 1121–1134 (2012).
10. Xie, M. *et al.* Age-related mutations associated with clonal hematopoietic expansion and malignancies. *Nat. Med.* **20**, 1472–1478 (2014).
11. Tomasetti, C., Marchionni, L., Nowak, M. A., Parmigiani, G. & Vogelstein, B. Only three driver gene mutations are required for the development of lung and colorectal cancers. *Proc. Natl. Acad. Sci.* **112**, 118–123 (2015).
12. Nordling, C. O. A new theory on the cancer-inducing mechanism. *Br. J. Cancer* **7**, 68–72 (1953).
13. Armitage, P. & Doll, R. The age distribution of cancer and a multi-stage theory of carcinogenesis. *Br. J. Cancer* **8**, 1–12 (1954).
14. Hyman, D. M. *et al.* The efficacy of larotrectinib (LOXO-101), a selective tropomyosin receptor kinase (TRK) inhibitor, in adult and pediatric TRK fusion cancers. *J. Clin. Oncol.* **35**, LBA2501-LBA2501 (2017).
15. Dees, N. D. *et al.* MuSiC: Identifying mutational significance in cancer genomes. *Genome Res.* **22**, 1589–1598 (2012).
16. Lawrence, M. S. *et al.* Mutational heterogeneity in cancer and the search for new cancer-associated genes. *Nature* **499**, 214–218 (2013).
17. Ciriello, G., Cerami, E., Sander, C. & Schultz, N. Mutual exclusivity analysis identifies oncogenic network modules. *Genome Res.* **22**, 398–406 (2012).
18. Leiserson, M. D. M. *et al.* Pan-cancer network analysis identifies combinations of rare somatic mutations across pathways and protein complexes. *Nat. Genet.* **47**, 106–114 (2014).
19. Paull, E. O. *et al.* Discovering causal pathways linking genomic events to transcriptional states using Tied Diffusion Through Interacting Events (TieDIE). *Bioinformatics* **29**, 2757–2764 (2013).
20. Bashashati, A. *et al.* DriverNet: uncovering the impact of somatic driver mutations on transcriptional networks in cancer. *Genome Biol.* **13**, R124 (2012).
21. Ng, S. *et al.* PARADIGM-SHIFT predicts the function of mutations in multiple cancers using pathway impact analysis. *Bioinformatics* **28**, 640–646 (2012).
22. Cheng, F., Zhao, J. & Zhao, Z. Advances in computational approaches for prioritizing driver mutations and significantly mutated genes in cancer genomes. *Brief. Bioinform.* **17**, 642–656 (2016).
23. Tokheim, C. J., Papadopoulos, N., Kinzler, K. W., Vogelstein, B. & Karchin, R. Evaluating the evaluation of cancer driver genes. *Proc. Natl. Acad. Sci.* **113**, 14330–14335 (2016).
24. Garraway, L. A. & Lander, E. S. Lessons from the cancer genome. *Cell* **153**, 17–37 (2013).
25. Kroschinsky, F. *et al.* New drugs , new toxicities : severe side effects of modern targeted and immunotherapy of cancer and their management. *Crit. Care* **21**, 1–11 (2017).
26. Park, S. R., Davis, M., Doroshow, J. H. & Kummar, S. Safety and feasibility of targeted agent combinations in solid tumours. *Nat. Rev. Clin. Oncol.* **10**, 154–168 (2013).
27. Hou, J. P. & Ma, J. DawnRank: Discovering personalized driver genes in cancer. *Genome Med.* **6**, 1–16 (2014).
28. Guo, W. F. *et al.* Discovering personalized driver mutation profiles of single samples in cancer by network control

- strategy. *Bioinformatics* **34**, 1893–1903 (2018).
29. Shrestha, R., Hodzic, E., Sauerwald, T. & Dao, P. HIT 'nDRIVE : Patient-Specific Multi-Driver Gene Prioritization for Precision Oncology. *Genome R* **27**, 1573–1588 (2017).
 30. Gitter, A. *et al.* Sharing information to reconstruct patient specific pathways in heterogeneous diseases. **8**, 1385–1395 (2014).
 31. Bailly-Bechet, M. *et al.* Finding undetected protein associations in cell signaling by belief propagation. *Proc. Natl. Acad. Sci. U. S. A.* **108**, 882–7 (2011).
 32. Szklarczyk, D. *et al.* STRING v10: Protein-protein interaction networks, integrated over the tree of life. *Nucleic Acids Res.* **43**, D447–D452 (2015).
 33. Joshi-Tope, G. *et al.* Reactome: A knowledgebase of biological pathways. *Nucleic Acids Res.* **33**, 428–432 (2005).
 34. Schaefer, C. F. *et al.* PID: The pathway interaction database. *Nucleic Acids Res.* **37**, 674–679 (2009).
 35. Ogata, H. *et al.* KEGG: Kyoto encyclopedia of genes and genomes. *Nucleic Acids Res.* **27**, 29–34 (1999).
 36. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol.* **15**, 1–21 (2014).
 37. McLachlen, G. & Peel, D. *Finite Mixture Models*. (Wiley inter-science, 2000).
 38. Forbes, S. A. *et al.* COSMIC: Somatic cancer genetics at high-resolution. *Nucleic Acids Res.* **45**, D777–D783 (2017).
 39. Koboldt, D. C. *et al.* Comprehensive molecular portraits of human breast tumours. *Nature* **490**, 61–70 (2012).
 40. Weinstein, J. N. *et al.* Comprehensive molecular characterization of urothelial bladder carcinoma. *Nature* **507**, 315–322 (2014).
 41. Muzny, D. M. *et al.* Comprehensive molecular characterization of human colon and rectal cancer. *Nature* **487**, 330–337 (2012).
 42. Lawrence, M. S. *et al.* Comprehensive genomic characterization of head and neck squamous cell carcinomas. *Nature* **517**, 576–582 (2015).
 43. Collisson, E. A. *et al.* Comprehensive molecular profiling of lung adenocarcinoma: The cancer genome atlas research network. *Nature* **511**, 543–550 (2014).
 44. Jonsson, P. F. & Bates, P. A. Global topological features of cancer proteins in the human interactome. *Bioinformatics* **22**, 2291–2297 (2006).
 45. Porta-Pardo, E., Garcia-Alonso, L., Hrabe, T., Dopazo, J. & Godzik, A. A Pan-Cancer Catalogue of Cancer Driver Protein Interaction Interfaces. *PLoS Comput. Biol.* **11**, 1–18 (2015).
 46. Dorsam, R. T. & Gutkind, J. S. G-protein-coupled receptors and cancer. *Nat. Rev. Cancer* **7**, 79–94 (2007).
 47. Cotto, K. C. *et al.* DGIdb 3.0: a redesign and expansion of the drug–gene interaction database. *Nucleic Acids Res.* **46**, 1068–1073 (2017).
 48. Van Allen, E. M. *et al.* Whole-exome sequencing and clinical interpretation of formalin-fixed, paraffin-embedded tumor samples to guide precision cancer medicine. *Nat. Med.* **20**, 682–688 (2014).
 49. Akhmedov, M. *et al.* PCSF : An R-package for network-based interpretation of high-throughput data. *PLoS Comput. Biol.* **13**, 1–7 (2017).