

RESEARCH ARTICLE

Open Access



Application of Random Forest and data integration identifies three dysregulated genes and enrichment of Central Carbon Metabolism pathway in Oral Cancer

Srija Mukhopadhyay¹, Sahana Ghosh¹, Debodipta Das¹, P. Arun², Bidyut Roy³, Nidhan K. Biswas¹, Arindam Maitra¹ and Partha P. Majumder^{1,3*} 

Abstract

Background: Studies of epigenomic alterations associated with diseases primarily focus on methylation profiles of promoter regions of genes, but not of other genomic regions. In our past work (Das et al. 2019) on patients suffering from gingivo-buccal oral cancer – the most prevalent form of cancer among males in India – we have also focused on promoter methylation changes and resultant impact on transcription profiles. Here, we have investigated alterations in non-promoter (gene-body) methylation profiles and have carried out an integrative analysis of gene-body methylation and transcriptomic data of oral cancer patients.

Methods: Tumor and adjacent normal tissue samples were collected from 40 patients. Data on methylation in the non-promoter (gene-body) regions of genes and transcriptome profiles were generated and analyzed. Because of high dimensionality and highly correlated nature of these data, we have used Random Forest (RF) and other data-analytical methods.

Results: Integrative analysis of non-promoter methylation and transcriptome data revealed significant methylation-driven alterations in some genes that also significantly impact on their transcription levels. These changes result in enrichment of the Central Carbon Metabolism (CCM) pathway, primarily by dysregulation of (a) *NTRK3*, which plays a dual role as an oncogene and a tumor suppressor; (b) *SLC7A5 (LAT1)* which is a transporter dedicated to essential amino acids, and is overexpressed in cancer cells to meet the increased demand for nutrients that include glucose and essential amino acids; and, (c) *EGFR* which has been earlier implicated in progression, recurrence, and stemness of oral cancer, but we provide evidence of epigenetic impact on overexpression of this gene for the first time.

(Continued on next page)

* Correspondence: ppm1@nibmg.ac.in

¹National Institute of Biomedical Genomics, Kalyani 741251, India

³Indian Statistical Institute, Kolkata, India

Full list of author information is available at the end of the article



© The Author(s). 2020 **Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>. The Creative Commons Public Domain Dedication waiver (<http://creativecommons.org/publicdomain/zero/1.0/>) applies to the data made available in this article, unless otherwise stated in a credit line to the data.

(Continued from previous page)

Conclusions: In rapidly dividing cancer cells, metabolic reprogramming from normal cells takes place to enable enhanced proliferation. Here, we have identified that among oral cancer patients, genes in the CCM pathway – that plays a fundamental role in metabolic reprogramming – are significantly dysregulated because of perturbation of methylation in non-promoter regions of the genome. This result compliments our previous result that perturbation of promoter methylation results in significant changes in key genes that regulate the feedback process of DNA methylation for the maintenance of normal cell division.

Keywords: Random Forest, Epigenomic, Transcriptomic, Integrative analysis, Gingivo-buccal oral cancer

Background

For various cancers both DNA methylation and gene expression data have been analyzed separately and alterations have been found to be associated with susceptibility and outcome [1, 2]. It is well known that DNA methylation impacts on gene expression. Therefore, attempts have been made to perform integrative analyses of these two types of data to draw robust inferences [3]. Various methods of data integration have been used [4, 5]. Methylation and expression data are high volume, highly correlated data. Further, the number of genes or DNA regions/sites on which data are collected are orders of magnitude higher than the number of patients and controls. This is commonly known as the “large p , small n ” problem or “curse of dimensionality” in Statistics. Many statistical methods involve inversion of a matrix for obtaining estimates of parameters. When the number of variables (p), on which data are available, for each patient exceeds the total number of patients (n), inversion of the relevant matrix becomes impossible [6]. This results in parameter estimates that are not unique; therefore, inferences are liable to be compromised. Random forest (RF) is a machine learning inferential method that is data-adaptive and tree-based. It handles correlated and large data sets very efficiently and is, therefore particularly appealing for analysis of high-dimensional genome data. Normally, only a small portion of a high-dimensional data is associated with a phenotype. A regression framework does not apply to this scenario. The highly correlated nature of genomic data also makes the application of standard statistical models inappropriate. RF is a non-parametric tree-based approach that is particularly suited for such data-analysis problems. RF can also be used to select and rank variables by taking advantage of variable importance measures. A good review of RF in genomic data analysis can be found [7].

We have used RF methodology to identify gene-body methylation differences between tumor and adjacent normal tissues in patients with oral squamous cell carcinoma of the gingivo-buccal region (OSCC-GB), the most common form of oral cancer in India [8, 9]. We then integrated the knowledge thus obtained with data on levels of transcription of genes, which we use as a

proxy for gene-expression levels, to discover methylation-driven alterations in the gene-body regions of the genome that significantly associate with dysregulation of genes in oral cancer.

DNA methylation occurs predominantly on cytosines followed by guanine residues (CpG). This type of methylation is referred to as CpG methylation. Although about 3–4% of all cytosines are methylated in normal human DNA, there are CpG islands, which are clusters of CpG dinucleotides in GC-rich regions, that remain unmethylated in all normal tissues [10]. Normally, a gene is transcribed if the CpG island in the promoter region remains unmethylated. But in cancer, the transcription of a tumor suppressor gene is silenced by the methylation of promoter CpG island of that gene. We had earlier analyzed data on methylation in CpG sites in the known promoter regions of all genes, but ignored gene-body CpG sites; sites that are on the coding regions of genes [4]. In our previous study, we identified about 200 genes that showed significant inverse correlation between promoter methylation and expression. These included a set of genes that act as transcription factors and genes associated with multiple cancer types. A significant finding of the study [4] was the identification of significant upregulation of CD274 and CD80 via promoter hypomethylation and hence immunosuppressive effects in OSCC-GB. Since in our previous study we had not considered gene body methylation, in the present study we have applied a modern data-adaptive method (RF) on gene-body methylation data and subsequently integrated with gene expression data. Our present analysis has resulted in the identification of some dysregulated genes and a pathway that were not identified in our earlier [4] analysis of promoter methylation and expression.

Methods

Patient recruitment and sample collection

This study was approved by the Institutional Ethics Committees of the Tata Medical Centre and the National Institute of Biomedical Genomics, India. Patients suffering from oral squamous cell carcinoma of the gingivobuccal region (OSCC-GB) were recruited into this

study with written informed consent. From each patient, a sample of tumor tissue and adjacent normal tissue were sampled by one of us (P.A.). The tissue samples were stored appropriately. TNM staging of 40 tumor samples were done following the 7th edition of the American Joint Committee on Cancer (AJCC) [11]. Summary statistics of demographic and clinical characteristics of the patients are provided in Table 1.

DNA methylation

Methylation data from paired tumour and adjacent normal tissue samples of 40 OSCC-GB patients were generated using the Illumina Infinium MethylationEPIC BeadChip [4]. Using the R package minfi, we estimated for each CpG site, the CpG-specific methylation level (β -value) as the ratio of the intensity of methylated (M) to the combined intensities of both methylated (M) and unmethylated (U) alleles:

$$\beta = \frac{M^*}{M^* + U^* + C}$$

where M^* and U^* denote signal intensities of M and U alleles, respectively, and the constant C set at 100 (as

recommended by the BeadChip manufacturer) [4, 12, 13]. The β -value ranges from 0 (unmethylated) to 1 (methylated). The sites that had a detection p -value ≥ 0.01 and those that mapped to X or Y chromosomes were removed. We further removed (a) probes that masked with “NA” values, (b) SNP associated probes with minor allele frequency (MAF) > 0.01 , (c) probes that overlapped with a repetitive element, (d) multi-mapped probes, (e) probes that did not map to annotated protein-coding genes [4, 12, 14, 15], and (f) probes that mapped to 3'UTR region of the genome.

Random Forest classifier

To analyze the difference between Tumor and Normal samples, a Random Forest (RF) method was used on Methylation data as implemented in the *randomForest* package in R [12, 14–19]. The random forest algorithm is an ensemble classifier similar to Classification and Regression Tree (CART) [17]. Each tree in an RF is built by choosing a bootstrap sample of two-third of the total number of individuals; the remaining one-third (Out-Of-Bag [OOB] sample) is utilised for validation. For each node in a tree, a binary splitting rule is used on a sample of CpG sites from the bootstrap sample to find the best split. The variable with the maximum information gain

Table 1 Demographic and clinical characteristics of 40 gingivo-buccal oral squamous cell carcinoma patients included in this study

Clinical Characteristics*	Frequency	Percent
Age		
< 40	7	0.18
40–50	16	0.40
51–60	11	0.28
> 60	6	0.15
Gender		
Male	33	0.83
Female	7	0.18
Risk-habit		
Chewing Tobacco	19	0.48
Chewing Tobacco and (Smoking and/or Alcohol)	16	0.40
Smoking and/or Alcohol	4	0.10
None	1	0.03
Tumor Stage		
T1	9	0.23
T2	12	0.30
T4	19	0.48
Lymph Node Invasion		
N0	21	0.53
N1	10	0.25
N2	9	0.23

*All patients were M0 (no metastasis) at the first presentation when tissue samples were collected for analysis

[20] is selected. A parameter *mtry* defines the number of variables randomly selected for each node in a tree, and another parameter *ntree* specifies the number of trees to be built in a forest. Normally, the value of *mtry* is taken to be the square root of the number of variables; this is also the default value in the R package. The output of *randomForest* provides an aggregated misclassification error (OOB error rate), which is estimated from predictions made on the OOB samples, and variable importance, which measures the weighted mean of the improvement in individual trees by each variable [15–17, 21]. The most reliable variable importance method is “permutation accuracy importance” or “Mean Decrease Accuracy” (MDA) [21, 22]. MDA permutes the data of *i*th variable in the OOB sample and records the permuted OOB error rate. The difference of the original and permuted OOB error rate averaged over the number of trees gives the importance score for *i*th variable (VI_i) in the random forest [19, 21–23]. A high value of MDA implies greater importance of the variable [21, 22].

$$VI_i = \frac{1}{ntree} \sum_{j=1}^{ntree} (OOBError_{ij}^{permuted} - OOBError_{ij})$$

Classification of samples

For efficient computation, only probes with $|\text{average } \Delta\beta| \geq 0.2$ were considered, where each $\Delta\beta$ was calculated by obtaining the difference between the β -values of tumor and adjacent normal samples of a patient for each probe indicating differential methylation between them and then taking average over the number of patients. A CpG site was considered hypermethylated if average $\Delta\beta \geq 0.2$ and hypomethylated if average $\Delta\beta \leq -0.2$ [4]. Before implementing the random forest (RF) classifier, *ntree* and *mtry* parameters were tuned to generate an accuracy rate [12, 16]. The best performing combination of parameters were those for which the OOB error rate stabilised and reached a minimum; i.e., the combination of parameters with the highest accuracy rate. Once the optimum set of parameters was determined, “randomForest” was executed 50 times on the methylation data of 40 paired samples. In each iteration variables (probes) with MDA-score > 0 were only selected [18]. The selected probes were then mapped to their respective genes. A gene was considered for further analyses if it satisfied the following conditions: (a) there were at least two probes in the non-promoter region of the gene, (b) methylation status of all probes in the non-promoter region were unidirectional; either hypermethylated or hypomethylated, and (c) had no probes in the promoter region. The stringency of criteria (a) and (b) were adopted to minimize the chance of false-positive

discovery, and the criterion (c) was adopted to make discoveries attributable to gene-body methylation only.

RNA sequencing

RNA was extracted and RNA sequencing was performed to obtain levels of transcription of genes, on the same set of 40 paired samples. Paired-end libraries were constructed and sequenced using Illumina HiSeq2500 [4, 24]. The quality of the RNA-Seq reads was checked by FastQC. *TopHat2* [4, 24–26] was then used to align these reads to a hg19 reference transcriptome or genome. Multi-mapped reads and non-concordant reads were filtered out using *SAMtools* [4] and duplicate reads were removed using *MarkDuplicates* from PICARD [4]. *Cufflinks* [4, 24–26] was then used to assemble and reconstruct the transcriptome. Finally, using *Cuffnorm*, normalised FPKM values for each gene were estimated [4]. Only those genes that had non-zero transcription levels in all samples were considered for further analysis. We have used the level of transcription of a gene as a proxy for the level of expression of the gene, and have used transcription and expression levels interchangeably in this report.

Integration of methylation and transcription data

Those genes for which there was no promoter probe and with multiple probes in the non-promoter region that were uniformly hyper- or hypo-methylated, and for which the level of transcription/expression change between tumour and normal tissues, averaged over the 40 pairs of samples, was higher than two-fold, were identified to be dysregulated by methylation in non-promoter regions [4]. Methylation effects on the 1st exon are similar to those of the promoter and exon boundary methylation modulates alternative splicing events [27]. Since this study is focused on gene expression alterations due to aberrant methylation on gene body, the genes that had 1st exon [28, 29] and exon boundary [27, 30] probes were removed. Finally, we considered only those genes for mapping on pathways that satisfied the known biological directionality of control; genes with hypermethylation (hypomethylation) in the gene-body region in the tumour tissue should have a significantly higher (lower) level of expression in the tumor tissue [31, 32].

Enrichment analysis of pathways

Genes that were so identified by the integration of both methylation and expression data were analyzed for enrichment of biological pathways. We considered pathways in KEGG for this analysis. ClueGo and CluePedia plug-ins of Cytoscape were used. To identify whether a pathway in KEGG was significantly enriched, a right-sided test based on hypergeometric distribution was used. Benjamini-Hochberg correction method was used to correct the *p*-values for multiple testing [4, 33].

Results

Identification of genes with abundant methylation in the non-promoter region

A total of 484,420 autosomal probes with detection p -value < 0.01 were associated with 18,688 genes. After removing 3'UTR and unannotated probes, 333,208 probes remained which were associated with 18,684 genes. Of these, 22,711 probes were with $|\text{average } \Delta\beta| \geq 0.2$ that mapped to 7027 genes. By fine-tuning (Figure S1), a stable OOB error rate was obtained with default $mtry = 150$ and $ntree = 2000$. Random forest was executed 50 times, with these optimal values of the parameters. The MDA scores of each variable and OOB error rate were recorded for 50 iterations. A uniform OOB error rate of 1.25% was observed in each iteration (Table S1). The set of probes with $MDA > 0$ comprised 10,105 probes that mapped to 4831 genes. Among these, for 433 genes all probes in the non-promoter region were hypermethylated, and for 233 genes all were hypomethylated. We have focused on these 666 unidirectionally methylated genes, for drawing further inferences integrated with gene expression patterns in tumor-normal paired tissues.

Integration of methylation and gene-expression

In paired tissues collected from the 40 OSCC-GB patients, non-zero levels of transcription/expression were found for 477 genes. Considering the 666 genes that exhibited significant and unidirectional methylation, it was found that 132 of these genes showed at least two-fold difference in the level of expression between tumour and normal tissues, averaged over the 40 patients. Of these 132 genes, 8 genes were removed as they had 1st exon and exon boundary probes. However, of these 124 only for 67 (54%) genes, the direction of change of expression level was consistent with that of methylation change (Table 2). That is, genes with hypermethylation (hypomethylation) in the tumour tissue had significantly higher (lower) levels of expression in the tumor tissue.

Enriched pathway

The pathway enrichment analysis using the 67 genes dysregulated by methylation alteration in the gene-body region between tumour and normal tissues, identified enrichment of one significant (corrected p -value = 0.0012) KEGG pathway. This was Central Carbon metabolism in Cancer with three associated genes EGFR, NTRK3, and SLC7A5. It has been reported, based on cell line studies, that overexpression of EGFR can impact on the development of solid tumors, including oral cancer [34]. It was found that EGFR was overexpressed and globally hypermethylated.

Discussion

By applying the novel Random Forest data-adaptive method to high-dimensional data (about 500,000 data points per individual) to identify significant alterations in gene-body methylation in gingivo-buccal oral tumor tissue compared to adjacent normal tissue, and subsequent integration with gene expression data it was detected that some genes and pathways were not earlier inferred to be involved in OSCC-GB only through cell-line studies. Although we found that only about 54% of genes found to have aberrant methylation were also dysregulated in the expected direction, this is not unexpected because gene-body methylation may not be the only cause of dysregulation of a gene. Hence, the directionality of dysregulation may not be in accord with what is expected under the methylation-transcription model. As a matter of fact, it is striking that over 50% of genes show transcription levels in accord with what is expected under gene-body hyper- or hypo-methylation. The significantly enriched pathway that has been identified using this data-adaptive and data-integrative approach is the Central Carbon Metabolism (CCM) pathway, which is involved in transport and oxidation of main carbon sources inside the cell. Fundamental cellular processes require energy for growth. The catabolic and anabolic reactions in metabolism are finely balanced and tightly regulated. Dysregulation results in cellular transformation and tumor progression. In rapidly dividing cancer cells, metabolic reprogramming from normal cells takes place to enable enhanced proliferation. CCM pathway plays a fundamental role in metabolic reprogramming. Changes in central carbon metabolism of cancer stem cells have also been noted [35]. It is noteworthy that enrichment of the CCM pathway in OSCC-GB takes place by gene-body methylation mediated dysregulation of three key genes, EGFR, NTRK3, SLC7A5 (Fig. 1).

Significant downregulation of NTRK3 mediated by promoter methylation was noted in our earlier study [4]. NTRK3 is a neurotrophin receptor. It behaves as an oncogene in breast cancer [36, 37] and possibly also in hepatocellular carcinoma [38]. However, it also plays a dual function. It acts as a tumor suppressor in colorectal cancer in which it is epigenetically inactivated [39]. In OSCC-GB also, NTRK3 is epigenetically dysregulated and appears to behave as a tumor suppressor.

SLC7A5 – earlier known as LAT1 – is a transporter dedicated to essential amino acids. Cancer cells have an increased demand for nutrients that include glucose and essential amino acids; the so-called “Warburg effect.” Overexpression of SLC7A5, as we have observed here, is explained in part by the presence, in its promoter, of a canonical binding site for the proto-oncogene *c-Myc* [40] that is known to regulate glucose metabolism [41].

Table 2 Results of 67 genes that showed significant relationship between methylation in the non-promoter region and gene expression

Gene	Mean of $\Delta\beta$ values of probes in the non-promoter region averaged over all patients	log2 fold-change of gene-expression values averaged over all patients
ABCA3	-0.242	-2.740
ADAMTS17	-0.235	-1.418
ADCY2	-0.282	-4.311
ADCYAP1R1	-0.231	-3.227
AFAP1L2	0.278	1.749
AGRN	0.299	1.830
ANGPT1	-0.255	-1.084
ANK2	-0.292	-3.425
ARNT2	-0.257	-1.056
ATP8A1	-0.263	-1.419
BCL11B	0.344	1.154
BMPER	-0.226	-2.030
BNC2	-0.284	-2.276
CACNA1D	-0.261	-2.461
CACNA2D1	-0.256	-2.616
CADM1	-0.256	-1.519
CDCA7	0.235	1.265
CIT	0.274	1.201
CLIC5	-0.264	-3.547
COBL	-0.283	-3.073
COL27A1	0.323	1.918
DNAH17	0.228	3.230
EEPD1	-0.245	-1.291
EGFR	0.240	1.147
EPHB2	0.285	2.239
EPSTI1	0.269	2.851
EXT1	0.272	1.068
FAM13C	-0.340	-1.862
FAM171A1	-0.232	-1.712
FGD5	-0.266	-1.149
FHIT	-0.289	-2.022
GF11	0.318	1.787
ICAM5	0.301	1.476
IGDCC4	-0.296	-2.178
KCNAB1	-0.279	-1.560
LAMB4	-0.206	-1.731
LDB2	-0.268	-1.250
LRP8	0.227	1.393
MCF2L	-0.242	-2.238
MEGF11	-0.257	-1.094

Table 2 Results of 67 genes that showed significant relationship between methylation in the non-promoter region and gene expression (*Continued*)

Gene	Mean of $\Delta\beta$ values of probes in the non-promoter region averaged over all patients	log2 fold-change of gene-expression values averaged over all patients
NCS1	0.290	1.266
NDRG1	0.208	1.343
NKAIN1	-0.241	-1.384
NTRK3	-0.298	-3.202
PALM	-0.274	-2.500
PAPPA	0.258	1.324
PARK2	-0.278	-2.854
PDZRN3	-0.250	-1.409
PLCL1	-0.243	-1.626
PML	0.220	1.562
PPM1L	-0.303	-2.386
PRKD1	-0.276	-1.051
RGS20	0.221	2.582
RTKN	0.235	1.183
SCIN	-0.278	-3.784
SDK2	0.247	2.410
SLC6A17	-0.240	-1.586
SLC7A5	0.241	1.367
SOBP	-0.286	-2.687
SPRED3	0.280	1.989
SUSD4	-0.319	-1.780
TECTA	-0.286	-1.007
TENM2	0.275	2.951
TMEM232	-0.213	-1.995
TRAM2	0.264	1.327
WNK2	-0.217	-3.895
ZNF423	-0.254	-1.799

Overexpression of SLC7A5 is also controlled by methylation in the promoter [4] and non-promoter regions (this study).

EGFR has been earlier implicated in progression, recurrence, and stemness of oral cancer [42, 43]. EGFR is inappropriately activated in cancer mainly because of amplification and point mutations. Transcriptional up-regulation of EGFR due to autocrine/paracrine mechanisms has also been described [44]. Here, for the first time, it is shown that dysregulation of EGFR takes place by epigenetic mechanisms in oral cancer.

Cancer cells rapidly multiply. Significant metabolic changes occur during cancer development and progression. Cancer cells have a lot of metabolic requirement

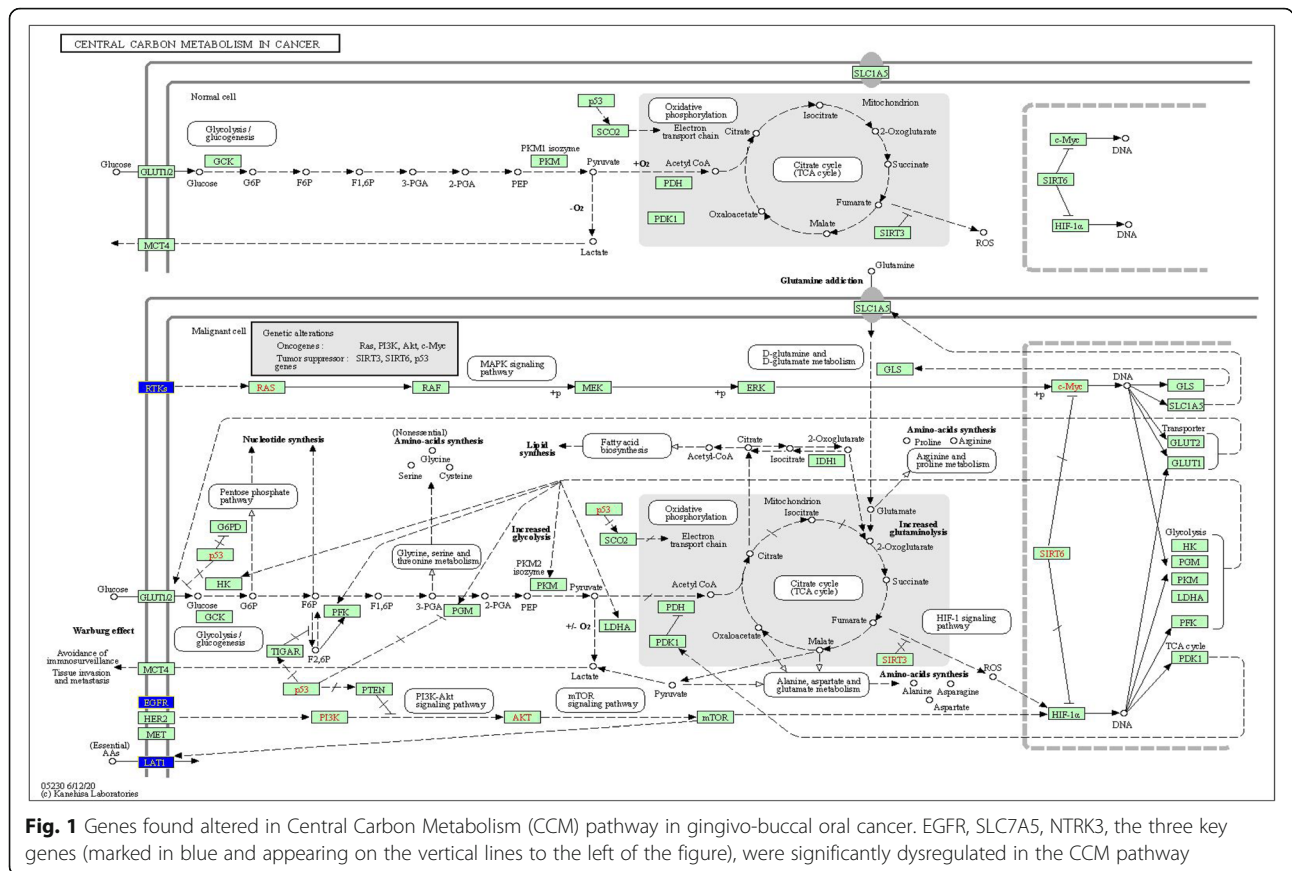


Fig. 1 Genes found altered in Central Carbon Metabolism (CCM) pathway in gingivo-buccal oral cancer. EGFR, SLC7A5, NTRK3, the three key genes (marked in blue and appearing on the vertical lines to the left of the figure), were significantly dysregulated in the CCM pathway

for the increase of their biomass and genome. These include the increased demand for nutrients such as glucose, essential amino acids and also glutamine, that becomes conditionally essential, for protein synthesis and/or energy supply [45–48]. Cancer cells utilize large amounts of glucose and glutamine and maintain high rates of glycolysis and glutaminolysis; called the Warburg effect. These increased requirements are met by the cancer cells themselves. Cancer cells undergo a large number of mutations, some of which take place in genes that belong to specific pathways, such as the central carbon metabolism pathway, which help meet these additional requirements. The central carbon metabolism pathway is large, complex and performs a variety of functions. About 70 genes are involved in this pathway, that are involved in a variety of functions that include glucose import, glycolysis, pentose phosphate flux, lactate excretion, pyruvate dehydrogenase flux, TCA cycle flux, pyruvate carboxylase flux, gluconeogenic flux, glycine biosynthesis, glutathione biosynthesis, proline biosynthesis, palmitate biosynthesis (fatty acid synthase activity), desaturation of palmitate, elongation of palmitate, and desaturation of stearate [49]. Changes in one or more components of the central carbon metabolism

pathway have been identified in various cancers. The genes that we have found to be significantly altered in their levels of expression resulting from gene-body methylation changes – NTRK3, SLC7A5 (LAT1) and EGFR – belong to the subcomponents related to the Warburg effect, notably glucose import and glycolysis.

Conclusions

Three key genes NTRK3, SLC7A5 (LAT1) and EGFR were dysregulated in the CCM pathway. Of these, NTRK3 [4, 50] and SLC7A5 [4] were earlier identified to be associated with oral cancer. However, we provide the first evidence of epigenetic impact on overexpression of EGFR in oral cancer. To enable enhanced proliferation of cells in a cancer tissue, metabolic reprogramming from normal cells usually takes place. In the present analysis, we have identified that among oral cancer patients, genes in the CCM pathway – that plays a fundamental role in metabolic reprogramming – are significantly dysregulated because of perturbation of methylation in non-promoter regions of the genome. This result compliments our previous result that perturbation of promoter methylation results in significant

changes in key genes that regulate the feedback process of DNA methylation for the maintenance of normal cell division. Taken together, it is evident that in oral cancer methylation driven alterations in both promoter and non-promoter genomic regions result in disruption of normal cell division accompanied by metabolic reprogramming to enable rapid cell proliferation.

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s12885-020-07709-0>.

Additional file 1: Figure S1. Fine-tuning *randomForest* parameters *mty* and *ntree*.

Additional file 2: Table S1. Results of 50 iterations of tuned Random Forest classifier showing the number of variables selected and the OOB error rate.

Abbreviations

RF: Random Forest; OSCC-GB: Oral Squamous Cell Carcinoma of the Gingivo Buccal region; CCM: Central Carbon Metabolism pathway; MAF: Minor Allele Frequency; CART: Classification And Regression Trees; OOB: Out-Of-Bag sample; MDA: Mean Decrease Accuracy; VI: Variable Importance; FPKM: Fragments Per Kilobase of transcript per Million mapped reads; KEGG: Kyoto Encyclopedia of Genes and Genomes

Acknowledgements

We are grateful to all participating members of Systems Medicine Cluster (SyMeC) and International Cancer Genome Consortium (ICGC) India project for their guidance and advice during the course of this study.

Authors' contributions

PPM and SM conceived of the study and designed the analysis of data. PA coordinated patient recruitment and sample collection. AM, NKB, DD and SG coordinated data generation and data collection. PPM and BR contributed towards the fine-tuning of the method. PPM and SM wrote the manuscript. PPM, DD, SG, PA, BR, AM and NKB edited the draft manuscript. All authors read and approved the final manuscript.

Funding

This work was partially supported by the J.C. Bose National Fellowship provided to PPM and a grant from the Department of Biotechnology (DBT), Govt. of India, through the SyMeC project (BT/Med-II/NIBM/G/SyMeC/2014/ Vol. II). The funding bodies played no role in the design of the study and collection, analysis, and interpretation of data and in writing the manuscript.

Availability of data and materials

Raw IDAT files of 40 samples generated using Illumina Infinium methylation array were deposited under EGAS00001003896 EGA study ID and aligned bam files for transcriptome data of 40 samples were deposited under EGAS00001003893 EGA study ID. Biospecimens may be shared on request, if not exhausted. Dispatch of biospecimens requires prior approval from the Government of India.

Ethics approval and consent to participate

The study was approved by the Institutional Ethics Committees, Dr. R. Ahmed Dental College & Hospital (RADCH), Kolkata, Chittaranjan National Cancer Institute (CNCI), Kolkata, National Institute of Biomedical Genomics (NIBM/G), Kalyani, and Indian Statistical Institute (ISI), Kolkata. Prior written informed consent was obtained from all study participants.

Consent for publication

Not applicable.

Competing interests

The authors declare that they have no competing interest.

Author details

¹National Institute of Biomedical Genomics, Kalyani 741251, India. ²Tata Medical Centre, Kolkata, India. ³Indian Statistical Institute, Kolkata, India.

Received: 26 July 2020 Accepted: 3 December 2020

Published online: 14 December 2020

References

- Jiao, Y., Widschwendter, M., Teschendorff, A.E. A systems-level integrative framework for genome-wide DNA methylation and gene expression data identifies differential gene expression modules under epigenetic control. *Bioinformatics*. 2014; 30:2360–2366. doi: <https://doi.org/10.1093/bioinformatics/btu316>.
- Udali S, Guarini P, Ruzzenente A, Ferrarini A, Guglielmi A, Lotto V, Tononi P, Pattini P, Moruzzi S, Campagnaro T, Conci S, Olivieri O, Corrocher R, Delledonne M, Choi S-W, Friso S. DNA methylation and gene expression profiles show novel regulatory pathways in hepatocellular carcinoma. *Clin Epigenet*. 2015;7:43. <https://doi.org/10.1186/s13148-015-0077-1>.
- Li M, Balch C, Montgomery JS, Jeong M, Chung JH, Yan P, Huang T, HM KS, Nephew KP. Integrated analysis of DNA methylation and gene expression reveals specific signaling pathways associated with platinum resistance in ovarian cancer. *BMC Med Genomics*. 2009;2:34. <https://doi.org/10.1186/1755-8794-2-34>.
- Das D, Ghosh S, Maitra A, Biswas NK, Panda CK, Roy B, Sarin R, Majumder PP. Epigenomic dysregulation-mediated alterations of key biological pathways and tumor immune evasion are hallmarks of gingivo-buccal oral cancer. *Clin Epigenet*. 2019;11(1):178. <https://doi.org/10.1186/s13148-019-0782-2>.
- Ma X, Liu Z, Zhang Z, Huang X, Tang W. Multiple network algorithm for epigenetic modules via the integration of genome-wide DNA methylation and gene expression data. *BMC Bioinformatics*. 2017;18:72. <https://doi.org/10.1186/s12859-017-1490-6>.
- Johnstone IM, Titterton DM. Statistical challenges of high-dimensional data. *Phil Trans R Soc A*. 2009;367:4237–53. <https://doi.org/10.1098/rsta.2009.0159>.
- Chen X, Ishwaran H. Random forests for genomic data analysis. *Genomics*. 2012;99:323–9. <https://doi.org/10.1016/j.jygeno.2012.04.003>.
- Muttagi SS, Patil BR, Godhi AS, Arora DK, Hallikerimath SR, Kale AD. Clinicopathological factors affecting lymph node yield in Indian patients with locally advanced squamous cell carcinoma of mandibular Gingivo-Buccal sulcus. *Indian J Cancer*. 2016;53:239–43. <https://doi.org/10.4103/0019-509X.197724>.
- Pathak KA, Gupta S, Talole S, Khanna V, Chaturvedi P, Deshpande MS, Pai PS, Chaukar DA, D'Cruz AK. Advanced squamous cell carcinoma of lower gingivobuccal complex: patterns of spread and failure. *Head Neck*. 2005;27:597–602. <https://doi.org/10.1002/hed.20195>.
- Esteller M. The necessity of a human epigenome project. *Carcinogenesis*. 2006;27(6):1121–5. <https://doi.org/10.1093/carcin/bgl033>.
- Sobin LH, Gospodarowicz MK, Wittekind C. TNM classification of malignant tumors. 7. Oxford: Wiley-Blackwell; 2011. p. 336.
- Zhang W, Spector TD, Deloukas P, Bell JT, Engelhardt BE. Predicting genome-wide DNA methylation using methylation marks, genomic position, and DNA regulatory elements. *Genome Biol*. 2015;16:14. <https://doi.org/10.1186/s13059-015-0581-9>.
- Ma X, Wang Y-W, Zhang MQ, Gazdar AF. DNA methylation data analysis and its application to cancer research. *Epigenomics*. 2013;5(3):301–16. <https://doi.org/10.2217/epi.13.26>.
- Everson TM, Lyons G, Zhang H, Soto-Ramírez N, Lockett GA, Patil VK, Merid SK, Söderhäll C, Melén E, Holloway JW, Arshad SH, Karmaus W. DNA methylation loci associated with atopy and high serum IgE: a genome-wide application of recursive Random Forest feature selection. *Genome Med*. 2015;7:89. <https://doi.org/10.1186/s13073-015-0213-8>.
- Naue J, Hoefsloot HJC, Mook ORF, Rijlaarsdam-Hoekstra L, van der Zwalm MCH, Henneman P, Kloosterman AD, Verschure PJ. Chronological age prediction based on DNA methylation: massive parallel sequencing and random forest regression. *Forensic Sci Int*. 2017;31:19–28. <https://doi.org/10.1016/j.fsigen.2017.07.015>.
- Houseman EA, Christensen BC, Yeh R-F, Marsit CJ, Karagas MR, Wrensch M, Nelson HH, Wiemels J, Zheng S, Wiencke JK, Kelsey KT. Model-based clustering of DNA methylation array data: a recursive-partitioning algorithm

- for high-dimensional data arising as a mixture of beta distributions. *BMC Bioinformatics*. 2018;9:365. <https://doi.org/10.1186/1471-2105-9-365>.
17. Christensen BC, Houseman EA, Marsit CJ, Zheng S, Wrensch MR, Wiemels JL, Nelson HH, Karagas MR, Padbury JF, Bueno R, Sugarbaker DJ, Yeh R-F, Wiencke JK, Kelsey KT. Aging and Environmental Exposures Alter Tissue-Specific DNA Methylation Dependent upon CpG Island Context. *Plos Genet*. 2009;5(8):e1000602. <https://doi.org/10.1371/journal.pgen.1000602>.
 18. Yang Y, Nephew K, Kim S. A novel k-mer mixture logistic regression for methylation susceptibility modeling of CpG dinucleotides in human gene promoters. *BMC Bioinformatics*. 2012;13:S15. <https://doi.org/10.1186/1471-2105-13-S3-S15>.
 19. Archer KJ, Kimes RV. Empirical characterization of random forest variable importance measures. *Comput Stat Data Anal*. 2008;52(4):2249–60. <https://doi.org/10.1016/j.csda.2007.08.015>.
 20. Deng H, Runger G. Gene selection with guided regularized random forest. *Pattern Recogn*. 2013;46:3483–9. <https://doi.org/10.1016/j.patcog.2013.05.018>.
 21. Strobl, C, Boulesteix, A, Kneib, T., Augustin, T., Zeileis, A Conditional variable importance for random forests *BMC Bioinformatics* 2008; 9:307. doi: <https://doi.org/10.1186/1471-2105-9-307>.
 22. Grömping U. Variable importance assessment in regression: linear regression versus random Forest. *Am Stat*. 2009;63(4):308–19. <https://doi.org/10.1198/tast.2009.08199>.
 23. Strobl C, Boulesteix A, Zeileis A, Hothorn T. Bias in random forest variable importance measures: illustrations, sources and a solution. *BMC Bioinformatics*. 2007;8:25. <https://doi.org/10.1186/1471-2105-8-25>.
 24. Ghosh, S., Chan C-KK. Analysis of RNA-Seq Data Using TopHat and Cufflinks. *Methods Mol Biol*. 2016; 1374:339–361. doi: https://doi.org/10.1007/978-1-4939-3167-5_18.
 25. Chu Y, Corey DR. RNA sequencing: platform selection, experimental design, and data interpretation. *Nucleic Acid Therapeutics*. 2012;22(4):271–4. <https://doi.org/10.1089/nat.2012.0367>.
 26. Trapnell C, Roberts A, Goff L, Pertea G, Kim D, Kelley DR, Pimentel H, Salzberg SL, Rinn JL, Pachter L. Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and cufflinks. *Nat Protoc*. 2012;7:562–78. <https://doi.org/10.1038/nprot.2012.016>.
 27. Gelfman S, Cohen N, Yearin A, Ast G. DNA-methylation effect on cotranscriptional splicing is dependent on GC architecture of the exon-intron structure. *Genome Res*. 2013;23:789–99. <https://doi.org/10.1101/gr.143503.112>.
 28. Løkk K, Modhukur V, Rajashekar B, Martens K, Magi R, Kolde R, Koltsina M, Nilsson TK, Viló J, Salumets A, Tonisson N. DNA methylome profiling of human tissues identifies global and tissue-specific methylation patterns. *Genome Biol*. 2014;15:3248. <https://doi.org/10.1186/gb-2014-15-4-r54>.
 29. Brenet F, Moh M, Funk P, Feierstein E, Viale AJ, Socci ND, Scandura JM. DNA Methylation of the First Exon Is Tightly Linked to Transcriptional Silencing. *Plos One*. 2011;6(1):e14524. <https://doi.org/10.1371/journal.pone.0014524>.
 30. Sun X, Tian Y, Wang J, Sun Z, Zhu Y. Genome-wide analysis reveals the association between alternative splicing and DNA methylation across human solid tumors. *BMC Med Genomics*. 2020;13(4). <https://doi.org/10.1186/s12920-019-0654-9>.
 31. Jjingo D, Conley AB, Yi SV, Lunnyk W, Jordan I. On the presence and role of human gene-body DNA methylation. *Oncotarget*. 2012;3:462–74. <https://doi.org/10.18632/oncotarget.497>.
 32. Yang X, Han H, Carvalho D, D, D., Lay, F., D., Jones, P., A., Liang G. Gene body methylation can Alter gene expression and is a therapeutic target in Cancer. *Cancer Cell* 2014; 26:1–14. doi: <https://doi.org/10.1016/j.ccr.2014.07.028>.
 33. Bindea G, Mlecnik B, Hackl H, Charoentong P, Tosolini M, Kirilovsky A, Fridman W-H, Pages F, Trajanoski Z, Galon J. ClueGO: a Cytoscape plug-in to decipher functionally grouped gene ontology and pathway annotation networks. *Bioinformatics*. 2009;25(8):1091–3. <https://doi.org/10.1093/bioinformatics/btp101>.
 34. Huang C-Y, Chan C-Y, Chou I-T, Lien C-H, Hung H-C, Lee M-F. Quercetin induces growth arrest through activation of FOXO1 transcription factor in EGFR-overexpressing oral cancer cells. *J Nutr Biochem*. 2013;24(9):1596–603. <https://doi.org/10.1016/j.jnutbio.2013.01.010>.
 35. Wong TL, Che N, Ma S. Reprogramming of central carbon metabolism in cancer stem cells. *Biochim Biophys Acta (BBA) - Mol Basis Dis*. 2017;1863:1728–38. <https://doi.org/10.1016/j.bbadis.2017.05.012>.
 36. Li, Z., Tognon, CE, Godinho, FJ, Yasaitis L, Hock H, Herschkowitz JI, Lannon CL, Cho E, Kim S-J, Bronson RT, Perou CM, Sorensen PH, Orkin SH. ETV6-NTRK3 Fusion Oncogene Initiates Breast Cancer from Committed Mammary Progenitors via Activation of AP1 Complex. *Cancer Cell*. 2007; 12:542–558. doi: <https://doi.org/10.1016/j.ccr.2007.11.012>.
 37. Tognon C, Knezevich SR, Huntsman D, Roskelley CD, Melnyk N, Mathers JA, Becker L, Carneiro F, MacPherson N, Horsman D, Poremba C, Sorensen PHB. Expression of the ETV6-NTRK3 gene fusion as a primary event in human secretory breast carcinoma. *Cancer Cell*. 2002;2:367–76. [https://doi.org/10.1016/S1535-6108\(02\)00180-0](https://doi.org/10.1016/S1535-6108(02)00180-0).
 38. Xiong D, Sheng Y, Ding S, Chen J, Tan X, Zeng T, Qin D, Zhu L, Huang A, Tang H. LINC00052 regulates the expression of NTRK3 by miR-128 and miR-485-3p to strengthen HCC cells invasion and migration. *Oncotarget*. 2016; 7(30):47593–608. <https://doi.org/10.18632/oncotarget.10250>.
 39. Luo Y, Kaz AM, Kanngurn S, Welsch P, Morris SM, Wang J, Lutterbaugh JD, Markowitz SD, Grady WM. NTRK3 Is a Potential Tumor Suppressor Gene Commonly Inactivated by Epigenetic Mechanisms in Colorectal Cancer. *Plos Genet*. 2013;9(7):e1003552. <https://doi.org/10.1371/journal.pgen.1003552>.
 40. Hayashi K, Jutabha P, Endou H, Anzai N. C-Myc is crucial for the expression of LAT1 in MIA Paca-2 human pancreatic cancer cells. *Oncol Rep*. 2012;28(3):862–6. <https://doi.org/10.3892/or.2012.1878>.
 41. Kim JW, Zeller KI, Wang Y, Jegga AG, Aronow BJ, O'Donnell KA, Dang CV. Evaluation of myc E-box phylogenetic footprints in glycolytic genes by chromatin immunoprecipitation assays. *Mol Cell Biol*. 2004;24:5923–36. <https://doi.org/10.1128/MCB.24.13.5923-5936.2004>.
 42. Mirza, Y., Ali, S., M., A., Awan, M., S., Idress, R., Naem, S., Zahid, N., Qadeer, U. Overexpression of EGFR in Oral premalignant lesions and OSCC and its impact on survival and recurrence. *Oncomedicine*. 2018; 3:28–36. doi: <https://doi.org/10.7150/oncm.22614>.
 43. Lv X-X, Zheng X-Y, Yu J-J, Ma H-R, Hua C, Gao R-T. EGFR enhances the stemness and progression of oral cancer through inhibiting autophagic degradation of SOX2. *Cancer Med*. 2019;00:1–10. <https://doi.org/10.1002/cam4.2772>.
 44. Wilson KJ, Mill C, Lambert S, Buchman J, Wilson TR, Hernandez-Gordillo V, Gallo RM, LMC A, Settleman J, Riese DJ II. EGFR ligands exhibit functional differences in models of paracrine and autocrine signaling. *Growth Factors*. 2012;30(2):107–16. <https://doi.org/10.3109/08977194.2011.649918>.
 45. Ganapathy V, Thangaraju M, Prasad PD. Nutrient transporters in cancer: relevance to Warburg hypothesis and beyond. *Pharmacol Ther*. 2009;121:29–40. <https://doi.org/10.1016/j.pharmthera.2008.09.005>.
 46. Heiden, M., G., V., Cantley, L., C., Thompson, C., B. Understanding the Warburg effect: the metabolic requirements of cell proliferation. *Science*. 2009; 324(5930):1029–1033. doi: <https://doi.org/10.1126/science.1160809>.
 47. Bhutia, Y., D., Ganapathy, V. Glutamine transporters in mammalian cells and their functions in physiology and cancer. *Biochim Biophys Acta (BBA) - Mol Cell Res* 2016; 1863(10):2531–2539. doi: <https://doi.org/10.1016/j.bbamcr.2015.12.017>.
 48. Scalise M, Pochini L, Galluccio M, Console L, Indiveri C. Glutamine Transport and Mitochondrial Metabolism in Cancer Cell Growth. *Front Oncol*. 2017; 7(306). <https://doi.org/10.3389/fonc.2017.00306>.
 49. Richardson AD, Yang C, Osterman A, Smith JW. Central carbon metabolism in the progression of mammary carcinoma. *Breast Cancer Res Treat*. 2008; 110:297–307. <https://doi.org/10.1007/s10549-007-9732-3>.
 50. Campbell PJ, Getz G, Korbel JO, et al. Pan-cancer analysis of whole genomes. *Nature*. 2020;578:82–93. <https://doi.org/10.1038/s41586-020-1969-6>.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

Ready to submit your research? Choose BMC and benefit from:

- fast, convenient online submission
- thorough peer review by experienced researchers in your field
- rapid publication on acceptance
- support for research data, including large and complex data types
- gold Open Access which fosters wider collaboration and increased citations
- maximum visibility for your research: over 100M website views per year

At BMC, research is always in progress.

Learn more biomedcentral.com/submissions

