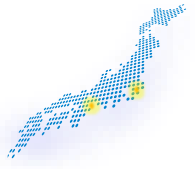


**FANTOM5**  
FUNCTIONAL ANNOTATION OF THE MAMMALIAN GENOME



# CAGE as a tool for cancer research and biomarker discovery

CAGE symposium

"New Topics of cancer research with CAGE method"

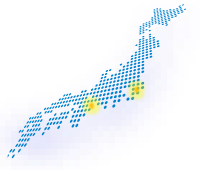
Bogumil Kaczkowski

Foreign Postdoctoral Researcher

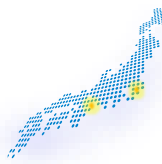
RIKEN CLST, Yokohama, Japan

Sept 13, 2016

# Agenda

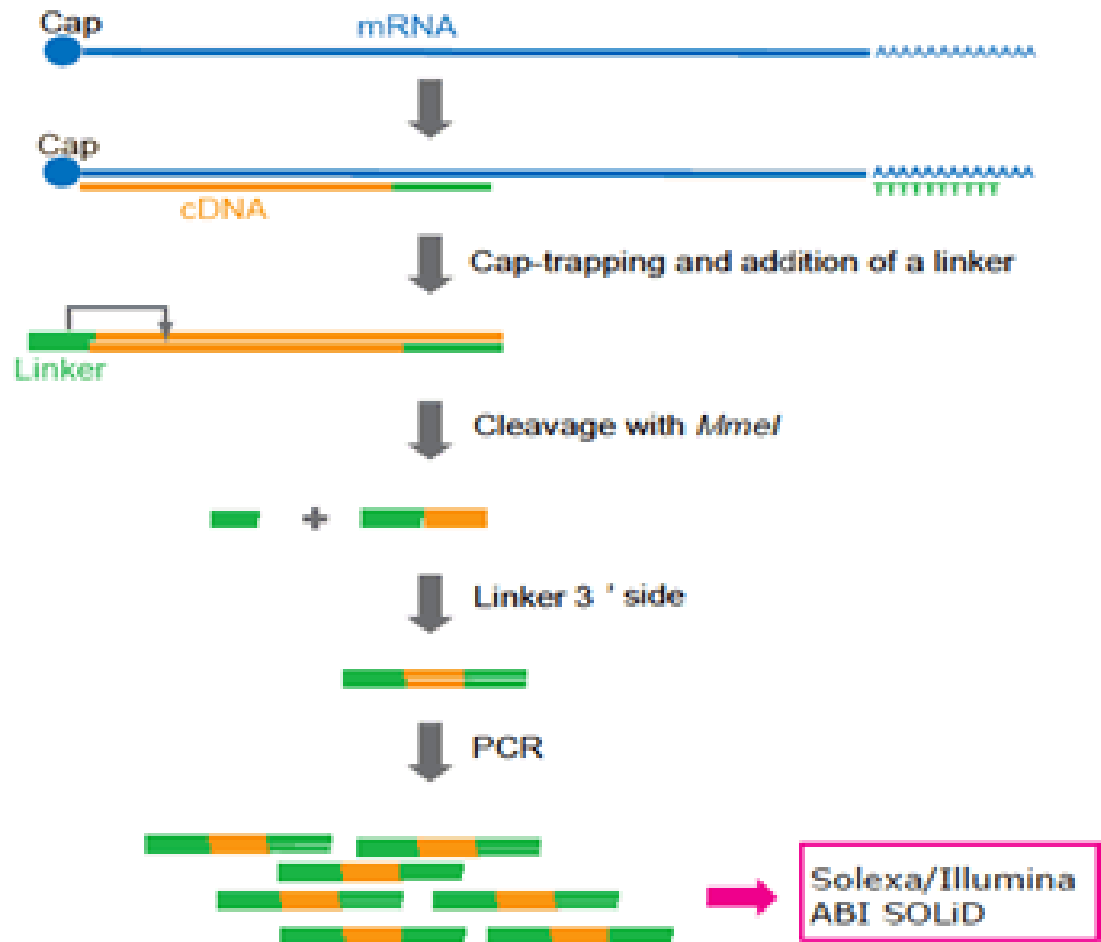


- Introduction to CAGE
- CAGE and biomarker discovery
- CAGE and cancer lncRNAs
- CAGE and enhancerRNAs
- CAGE and repetitive elements

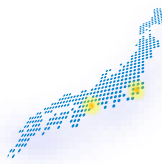


# Quick intro: Cap Analysis Gene Expression (CAGE)

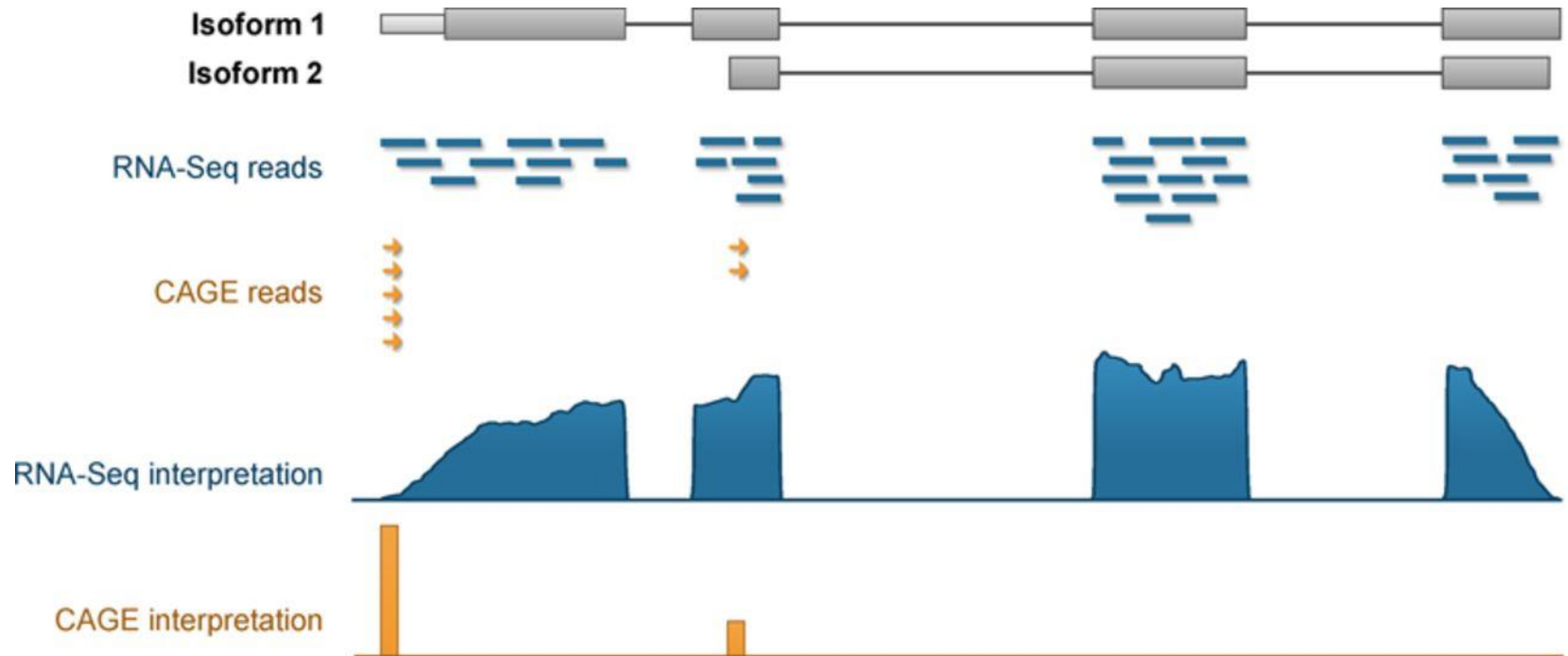
- Capture of 5' end of capped RNAs as a short sequence tags
- followed by high-throughput sequencing



<http://fantom.gsc.riken.jp/protocols/basic.html>

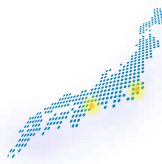


# CAGE Cap Analysis Gene Expression vs RNA-Sequencing

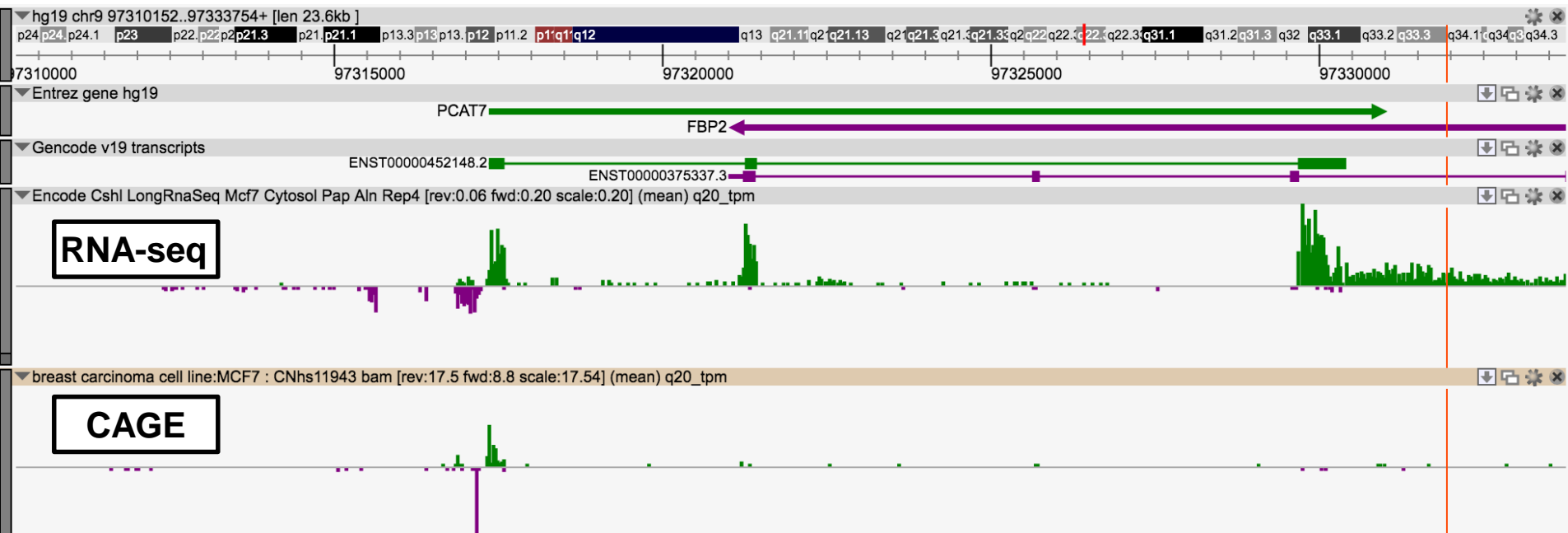


Nancy Yiu-Lin Yu et al. Nucl. Acids Res. 2015;nar.gkv608

Nucleic Acids Research



# CAGE Cap Analysis Gene Expression vs RNA-Sequencing a genome browser example

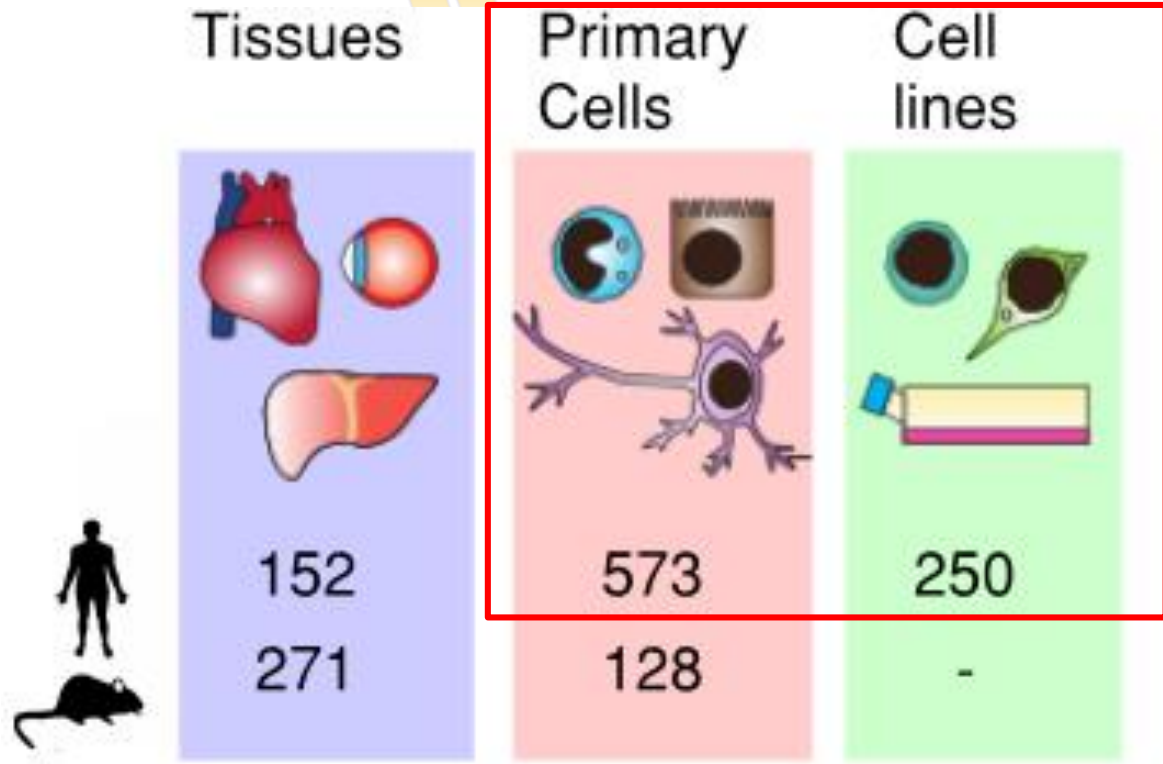


ZENBU genome browser  
<http://fantom.gsc.riken.jp/zenbu/>

Severin, J. et al. (2014). Interactive visualization and analysis of large-scale sequencing datasets using ZENBU. *Nature Biotechnology*, 32(3), 217–219.

RIKEN Center for Life Science Technologies





**Cancer vs. Normal Analysis**

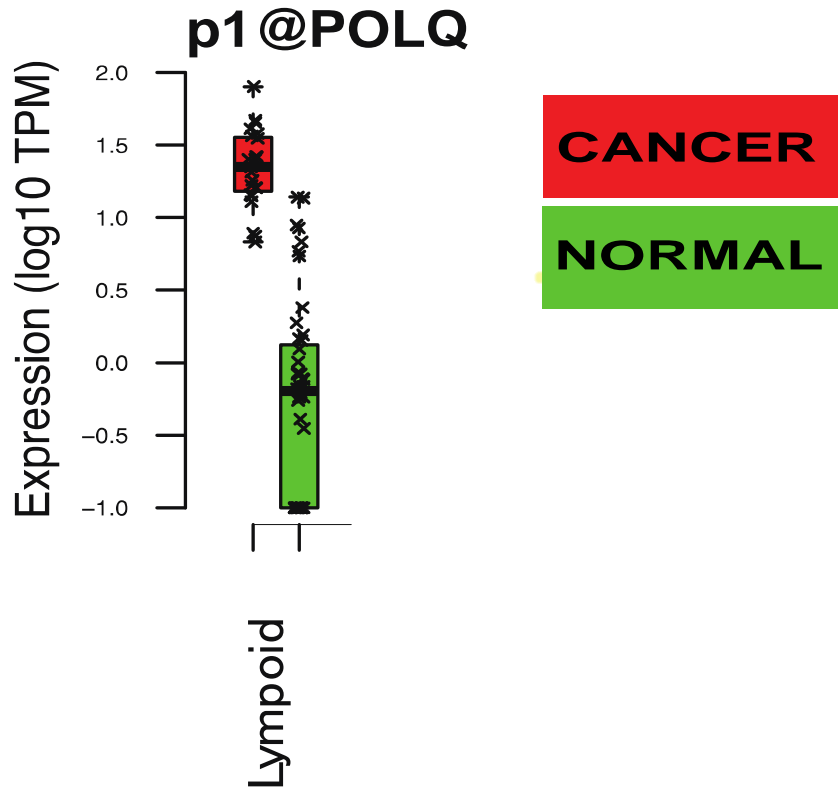
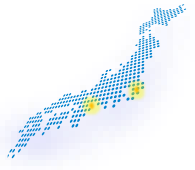
### Cap Analysis of Gene Expression

- CAGE is 5' sequence tag technology that globally determines transcription start sites (TSS) in the genome and their expression levels.
- CAGE was applied to a broad selection of cancer cell lines and primary cells

FANTOM Consortium and the RIKEN PMI and CLST (DGT). (2014).  
 A promoter-level mammalian expression atlas. *Nature*, 507(7493), 462–470.

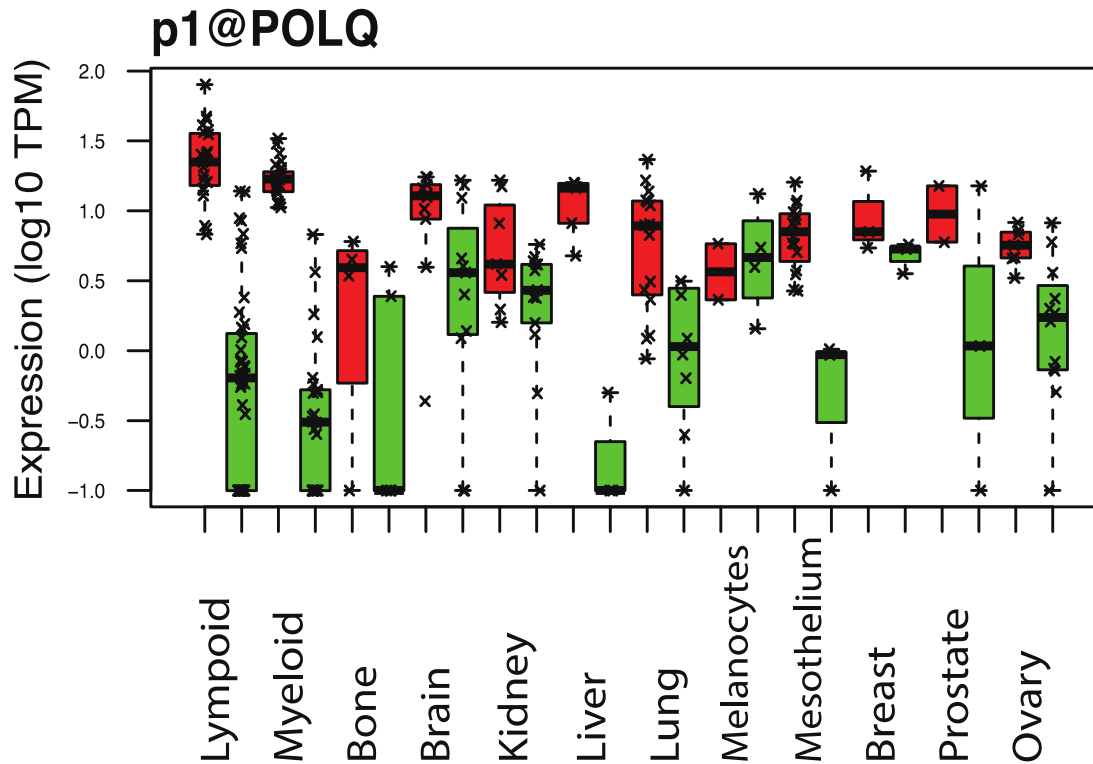
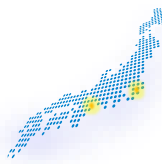


# Finding differentially expressed genes in cancer



- We compare **cancer cell lines** to **normal primary cells**
- But we also want to compare matching cancer-normal pairs

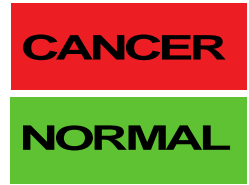
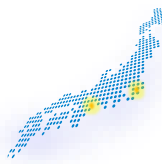
# Pan Cancer Differential Expression



- **Genewise Negative Binomial Generalized Linear Models** as implemented in **edgeR**
- Design matrix accounted for origin
- Contrast matrix to extract global cancer vs normal change
- FDR < 0.01
- Fold change > 2

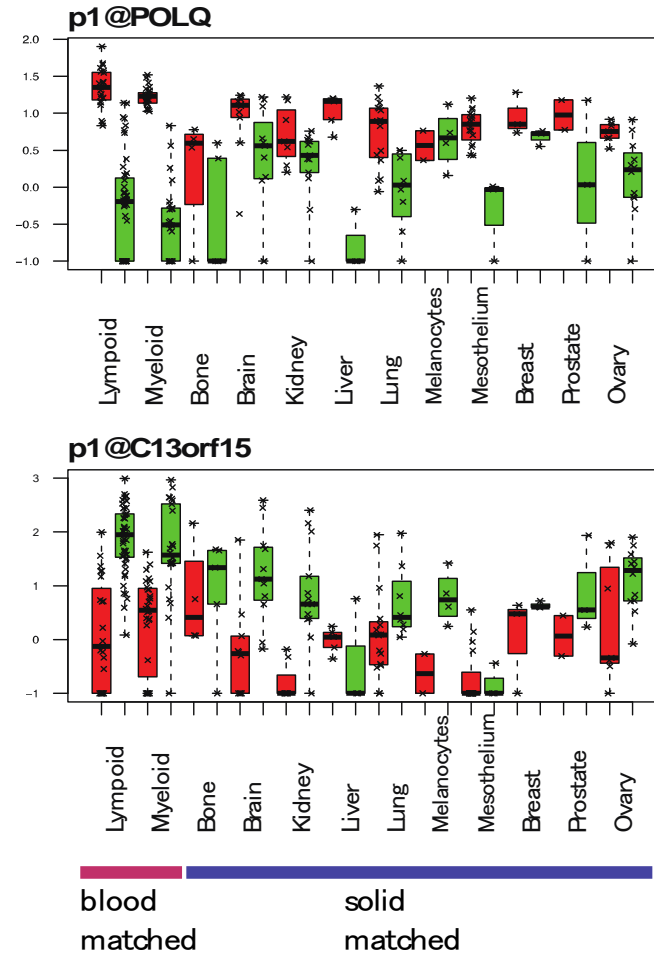
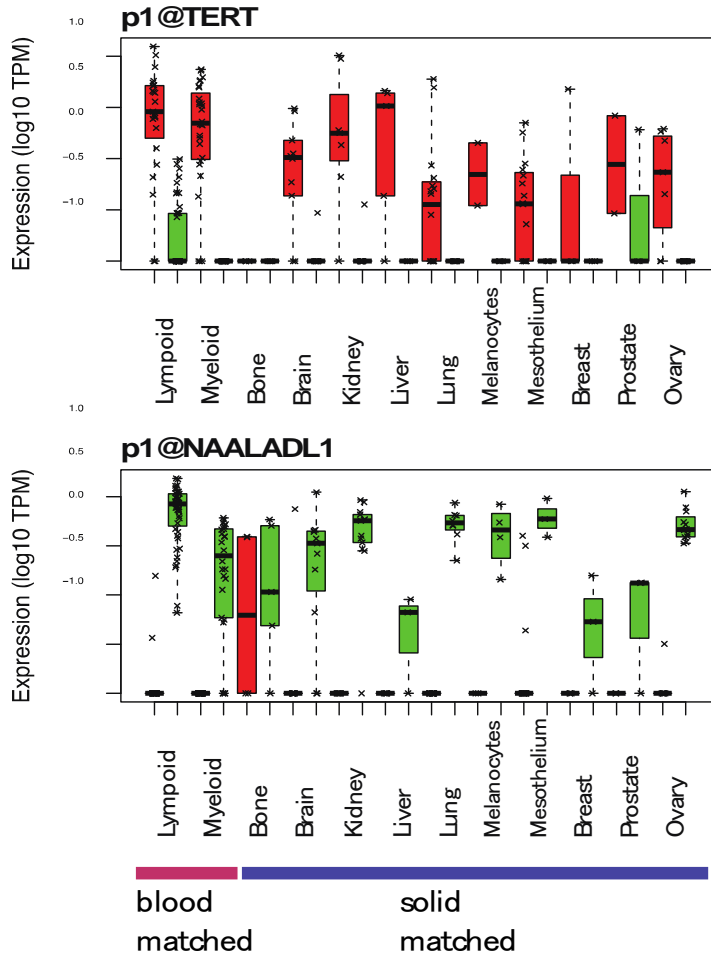


# Pan Cancer Differential Expression



**Switching ON/OFF**

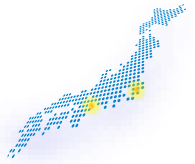
**Expression shift UP/DOWN**



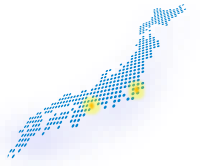
**Up regulated**

**Down regulated**

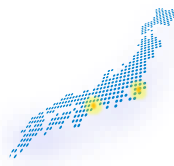
# Annotation of differentially expressed promoters



	UP-regulated		DOWN-regulated	
	# promoter	# genes	# promoter	# genes
protein coding	575	430	371	295
protein coding gene body	347	234	77	19
lncRNAs	209	148	18	14
antisense	65	52	4	2
pseudogene	34	29	6	6
small ncRNAs	14	13	0	0
not_annotated	351		37	
Sum	1595	906	513	336



We work with cancer cell lines,  
are the results relevant to  
clinical tumors?



# TCGA The Cancer Genome Atlas

- Available RNA-seq, gene expression profiles for:
  - ~ **5000 clinical tumor** samples
  - ~ 500 normal tissues
  - **14 cancer types**

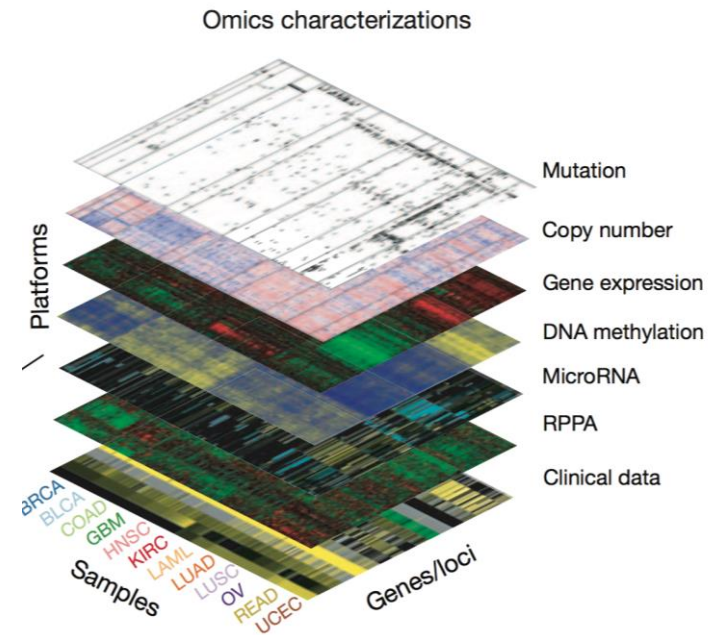
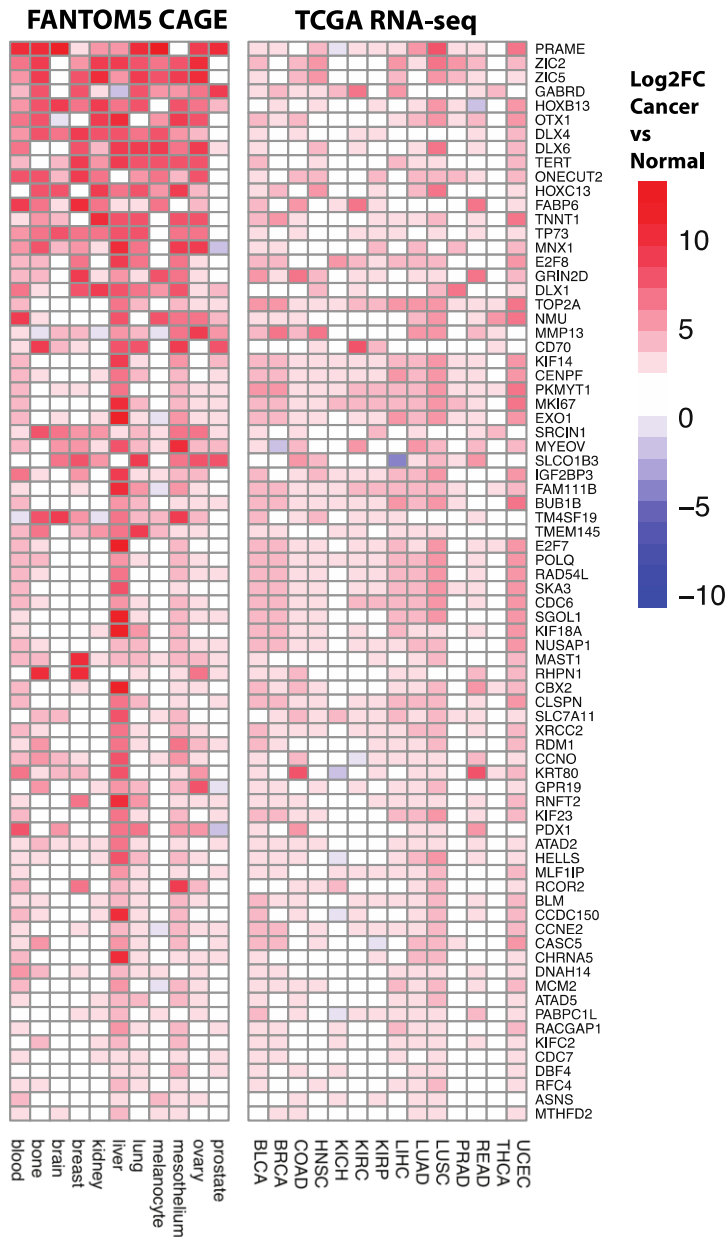
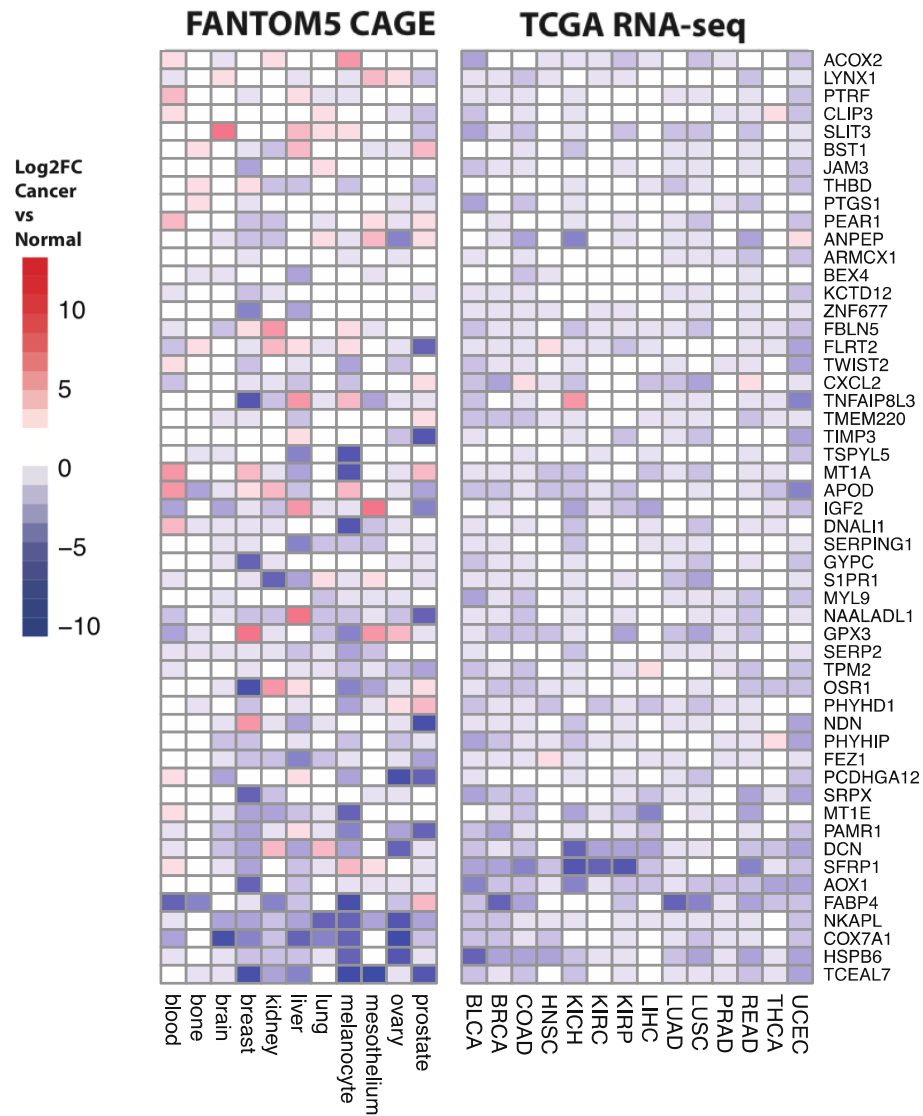
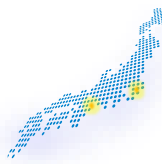


Figure 2



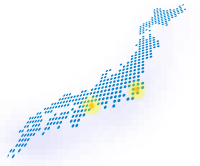
# Up regulated in both F5 and TCGA

- 76 genes (~17%)
- enriched in genes involved in:
  - cell cycle
  - DNA metabolism,
  - biopolymer metabolism
  - homeobox genes (developmental)



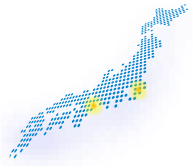
# Down regulated in both F5 and TCGA

- 54 genes (~20%)
- enriched in genes involved in:
  - oxidoreductase activity

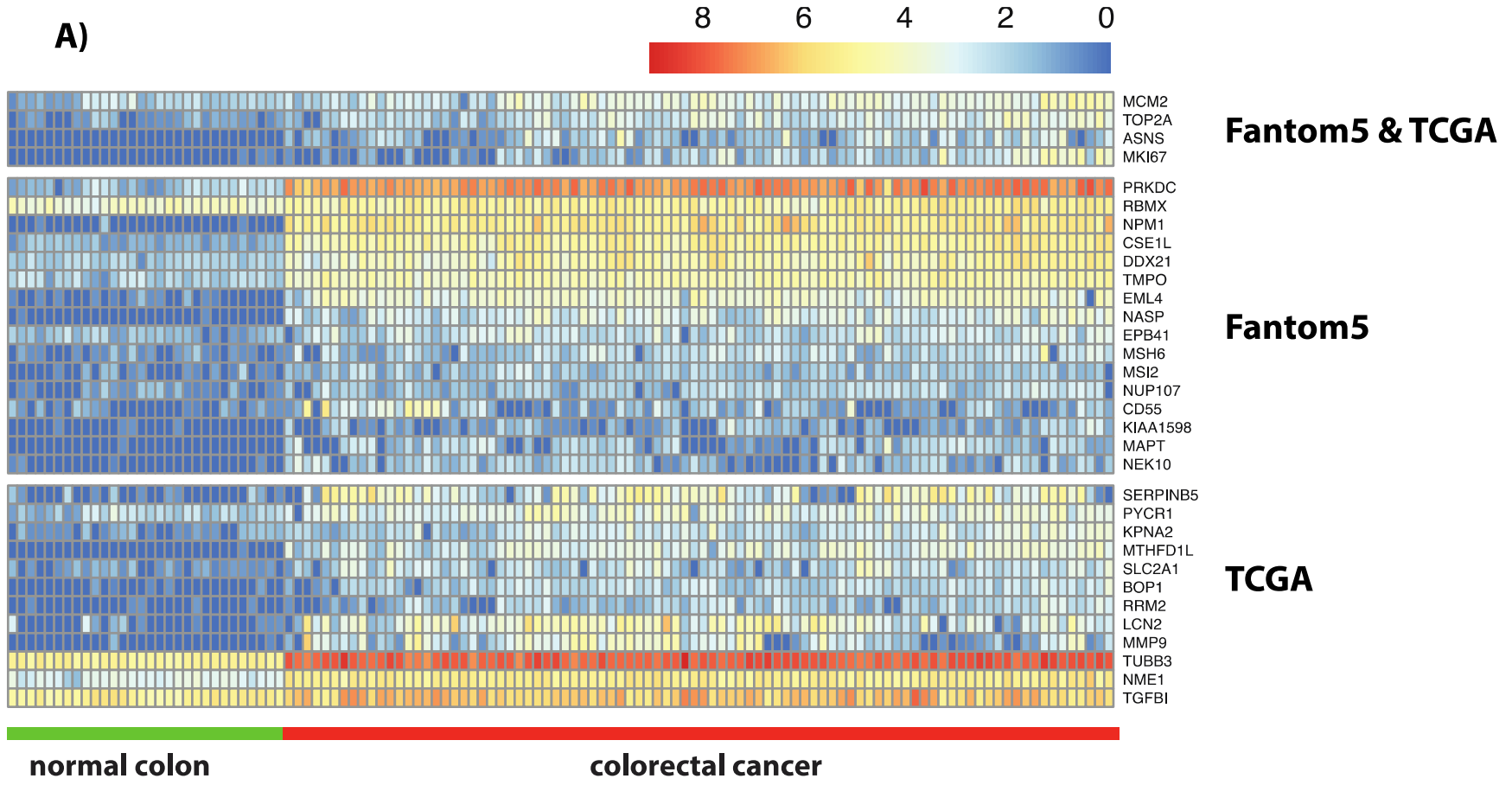


# Do RNA translate into protein?

Are the biomarker candidates de-regulated  
at protein level?

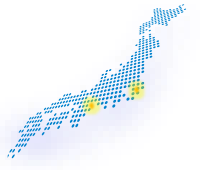


# Protein level conformation of cancer up-regulated PC genes



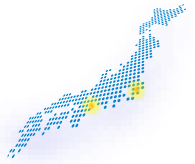
Spectral count data from 90 colorectal cancers and 30 normal (mass spec from Zhang *et al.* 2015)





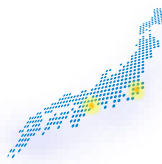
# Long non coding RNAs?

# Annotation of differentially expressed promoters

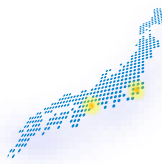


	UP-regulated		DOWN-regulated	
	# promoter	# genes	# promoter	# genes
protein coding	575	430	371	295
protein coding gene body	347	234	77	19
lncRNAs	209	148	18	14
antisense	65	52	4	2
pseudogene	34	29	6	6
small ncRNAs	14	13	0	0
not_annotated	351		37	
Sum	1595	906	513	336

# lncRNA confirmed in TCGA tumors (25 up, 3 down)

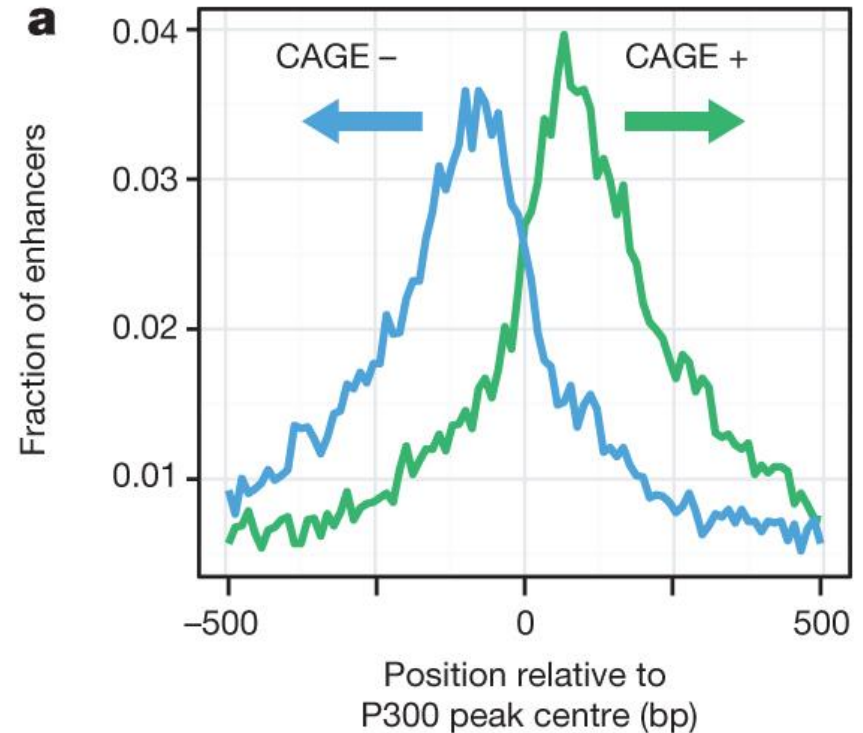


func_name miTranscriptome	gene_name (GENCODE19)	lncRNA type	mean tpm	FANTOM5 Pan-Cancer DE	# of TCGA cancers UP (TCGA)	# of TCGA cancers DOWN (TCGA)
	<b>PVT1</b>	lncrna	0.407	UP	9	0
CAT647	CTD-2023M8.1	linc	1.046	UP	5	0
CAT122	<b>RP11-284F21.7**</b>	antisense	3.201	UP	5	0
CAT2260	RP11-191L9.4	linc	0.404	UP	4	0
MNX1-AS1	MNX1-AS1	antisense	0.198	UP	4	0
CAT1138	G084254	linc	2.695	UP	3	0
CAT266	G044387	linc	0.103	UP	3	0
CAT62	RP4-792G4.2	antisense	0.178	UP	3	0
LINC00898	LINC00898	linc	0.104	UP	3	0
DGCR5	<i>DGCR5</i>	antisense	3.738	UP	3	0
CAT2039	AC009005.2	antisense	0.417	UP	3	0
CAT1022	G080198	linc	2.452	UP	3	0
CAT1167	<b>PCAT7</b>	antisense	0.086	UP	2	0
FEZF1-AS1	FEZF1-AS1	antisense	0.958	UP	2	0
CAT1572	RP11-438N16.1	linc	1.574	UP	2	0
DLX6-AS1	DLX6-AS2	antisense	0.259	UP	2	0
CAT1833	RP11-57A19.2	linc	0.528	UP	2	0
CAT615	G064032	linc	0.048	UP	2	0
LINC00669	LINC00669	linc	4.096	UP	2	0
CAT800	<b>RP11-328M4.2**</b>	antisense	0.287	UP	2	0
MF12-AS1	MF12-AS1	antisense	1.301	UP	2	0
CAT219	AC074117.10	antisense	1.447	UP	2	0
	LINC00511	linc	0.256	UP	3	0
	AC006262.5	linc	0.223	UP	2	0
	RP11-435O5.2	linc	0.475	UP	2	0
CAT1235	RP11-124N14.3*	antisense	7.352	DOWN	2	-4
MEG3	<b>MEG3</b>	linc	182.23	DOWN	0	-2
	MT1L	pseudogene	116.644	DOWN	0	-9

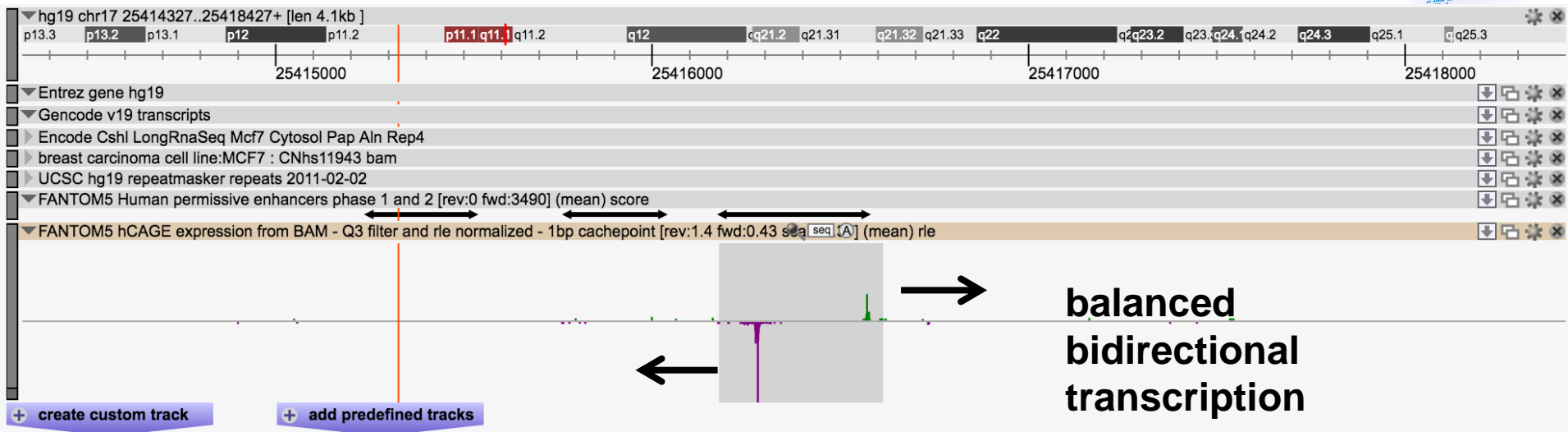
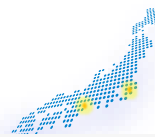


# Enhancer RNAs

- CAGE data can be used to estimate the activity of enhancers
- balanced bidirectional capped transcription)



# Example of "cancer enhancer"

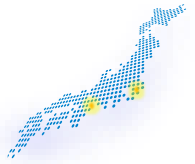


**FANTOM5 hCAGE expression from BAM - Q3 filter and rle normalized - 1bp cachepoint** chr17 25416178..25416609 (432bp) [STAT] [FLT] [ \* ] [ v ]

experiment name (988 / 988)

	<- anti-sense strand	> sense strand ->
chronic myelogenous leukemia cell line:K562 ENCODE, biol rep1 : CNhs12334 bam [rle]	205.52	87.98
small cell lung carcinoma cell line:WA-hT : CNhs11812 bam [rle]	346.91	79.26
chronic myelogenous leukemia cell line:K562 ENCODE, biol rep2 : CNhs12335 bam [rle]	149.6	66.47
chronic myelogenous leukemia cell line:K562 ENCODE, biol rep3 : CNhs12336 bam [rle]	143.49	63.1
signet ring carcinoma cell line:Kato III : CNhs10753 bam [rle]	177.46	55.26
fibrous histiocytoma cell line:GCT TIB-223 : CNhs11842 bam [rle]	164.34	45.9
gastric adenocarcinoma cell line:MKN1 : CNhs11737 bam [rle]	99.7	44.7
chondrosarcoma cell line:SW 1353 : CNhs11833 bam [rle]	380.52	44.2
neuroblastoma cell line:NH-12 : CNhs11811 bam [rle]	77.21	34.97
squamous cell carcinoma cell line:EC-GI-10 : CNhs11252 bam [rle]	184.98	32.98
pagetoid sarcoma cell line:Hs 925.T : CNhs11856 bam [rle]	154.1	31.97
fibrosarcoma cell line:HT-1080 : CNhs11860 bam [rle]	88.25	27.41
renal cell carcinoma cell line:OS-RC-2 : CNhs10729 bam [rle]	140.35	23.32
choriocarcinoma cell line:T3M-3 : CNhs11820 bam [rle]	65.94	22.6
non-small cell lung cancer cell line:NCI-H1385 : CNhs12193 bam [rle]	66.83	11.41
maxillary sinus tumor cell line:HSQ-89 : CNhs10732 bam [rle]	54.18	15.86
small cell cervical cancer cell line:HCSC-1 : CNhs11885 bam [rle]	127.03	18.67
squamous cell carcinoma cell line:T3M-5 : CNhs11739 bam [rle]	37.76	15.53
liposarcoma cell line:SW 872 : CNhs11851 bam [rle]	33.37	13.26
bronchioalveolar carcinoma cell line:NCI-H358 : CNhs11840 bam [rle]	70.43	13.22
lung adenocarcinoma cell line:PC-14 : CNhs10726 bam [rle]	93.04	11.98
bronchogenic carcinoma cell line:ChaGo-K-1 : CNhs11841 bam [rle]	38.53	11.58
breast carcinoma cell line:MCF7 : CNhs11943 bam [rle]	59.64	10.68
small cell lung carcinoma cell line:NCI-H82 : CNhs12809 bam [rle]	71.57	9.78
mucinous adenocarcinoma cell line:JHOM-1 : CNhs11752 bam [rle]	21.4	7.37
choriocarcinoma cell line:SCH : CNhs11875 bam [rle]	23.9	6.95
chronic myelogenous leukemia cell line:K562 : CNhs11250 bam [rle]	20.9	5.5
adult T-cell leukemia cell line:ATN-1 : CNhs10738 bam [rle]	14.15	3.68
oral squamous cell carcinoma cell line:Ca9-22 : CNhs10752 bam [rle]	8.24	2.94
meningioma cell line:HKBMM : CNhs11945 bam [rle]	5.65	2.82

expression/activation in cancer cells only

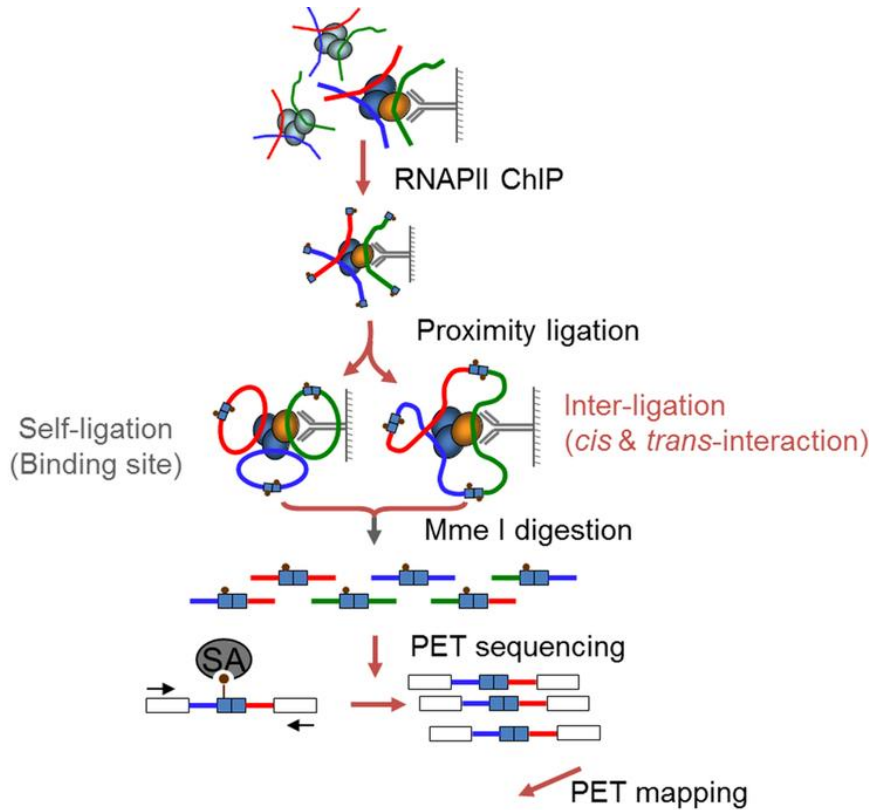


# Enhancer RNAs

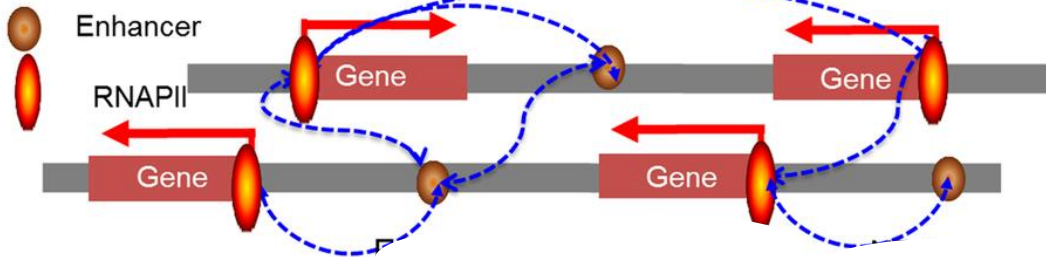
- differential expression analysis on CAGE tags counts under 43,011 enhancers
- We discovered 90 enhancer up regulated in cancer

?	Up-regulated?		Down-regulated?	
	Pan?Cancer?	Solid?Cancer?	Pan?Cancer?	Solid?Cancer?
Enhancers?	28?	62?	0?	0?

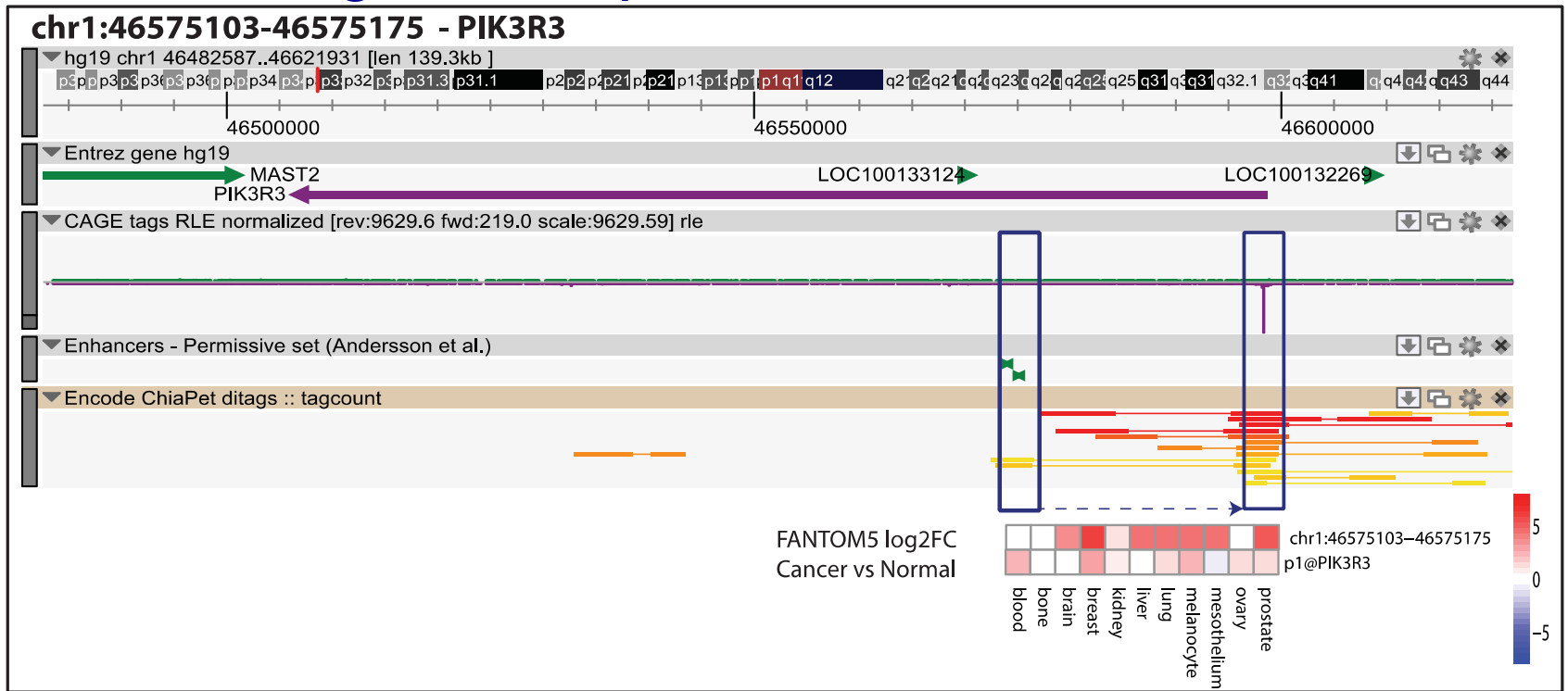
# Pan-cancer enhancers physically associated with genes implicated in cancer.



- Chia-PET data
- (Chromatin Interaction Analysis by Paired-End Tag Sequencing)
- from Encode project
- Physical chromatin associations of enhancers to promoters of genes related to cancer.



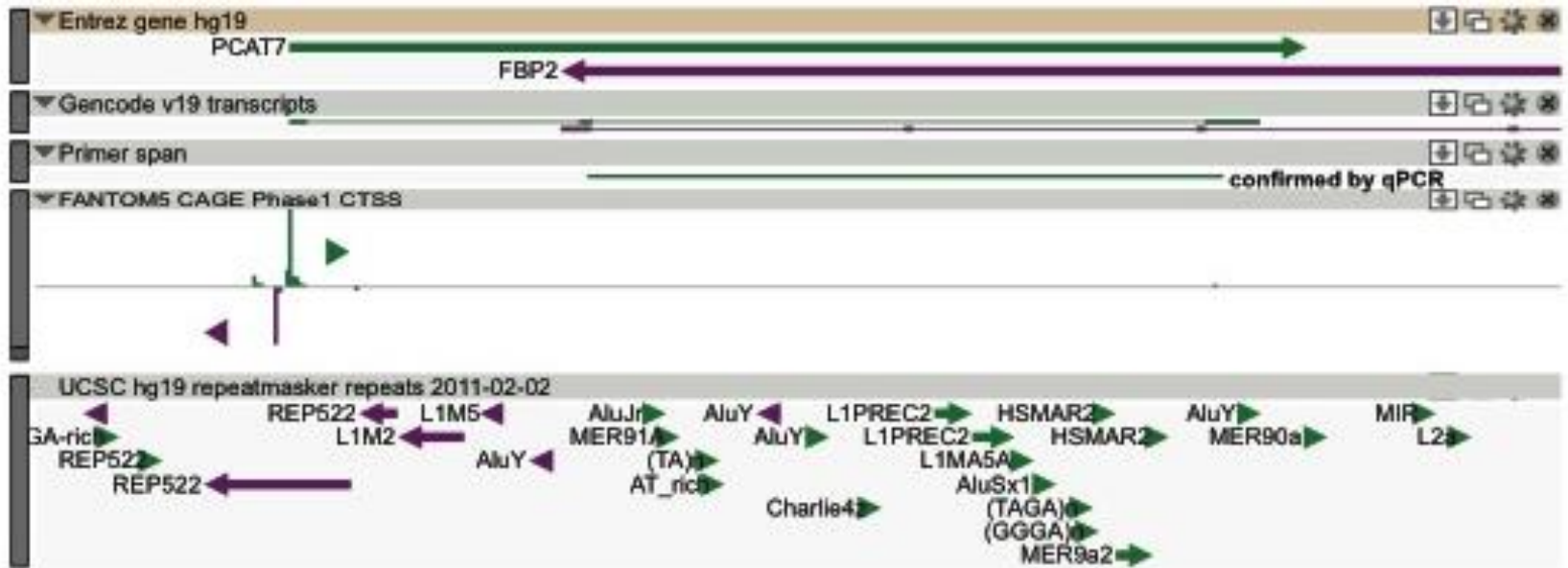
# Pan-cancer enhancers physically associated with genes implicated in cancer.



- Enhancer at chr1:46575103-46575175 is shown to physically associate with the promoter for the PIK3R3 gene.
- PIK3R3- phosphoinositide-3-kinase, regulatory subunit 3
- PIK3R3 was reported increase proliferation in colorectal, lung, gastric cancer, leukemia and glioma.

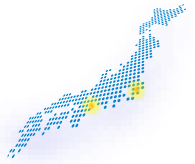


# CAGE and repetitive elements



**CAGE let's us see if promoter overlaps repeat elements**

# Promoters overlapping repetitive elements.



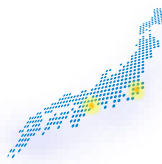
	# all overlapping promoters	Up-regulated			Protein coding						
		# promoters	Odd ratio	p-value	5' transcript	intronic	exon	3UTR	lncRNAs	Pseudogene	non annotated
<b>REP522</b>	72	<b>25</b>	<b>62.05</b>	<b>2.20E-16</b>	1	0	0	0	<b>9</b>	3	12
SINE/Alu	3961	<b>138</b>	<b>4.44</b>	<b>2.20E-16</b>	5	<b>67</b>	1	1	<b>11</b>	3	<b>50</b>
LTR/ERV1	3932	<b>133</b>	<b>4.30</b>	<b>2.20E-16</b>	7	<b>12</b>	0	0	<b>31</b>	2	<b>83</b>
LINE/L1	3426	<b>67</b>	<b>2.35</b>	<b>1.75E-09</b>	<b>2</b>	<b>22</b>	0	0	12	0	32
LTR/ERVL	1488	20	1.57	0.049	2	2	0	0	8	0	8
LINE/L2	3220	25	0.90	0.70	2	<b>4</b>	0	0	<b>4</b>	0	<b>17</b>
LTR/ERVL-MaLR	3613	31	0.99	1	6	<b>4</b>	0	0	<b>10</b>	0	<b>11</b>
Simple_repeat	11982	204	2.13	2.20E-16	86	70	4	7	17	1	63
Low_complexity	2013	18	1.04	0.81	<b>15</b>	2	2	6	2	0	4

Gencode v19 annotation

Kaczkowski et al. (2016). *Cancer Research*, 76(2), 216–

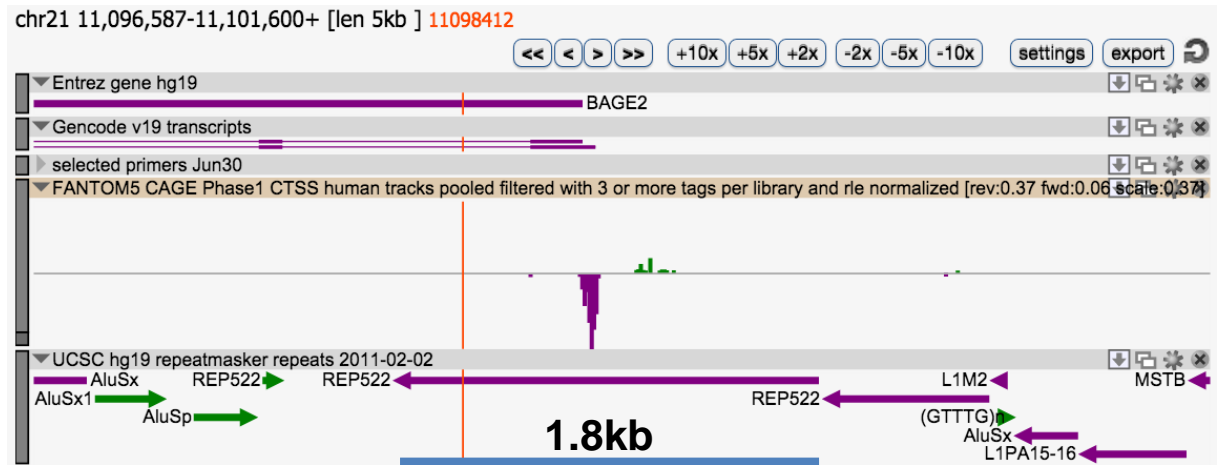
226.

# REP522 story

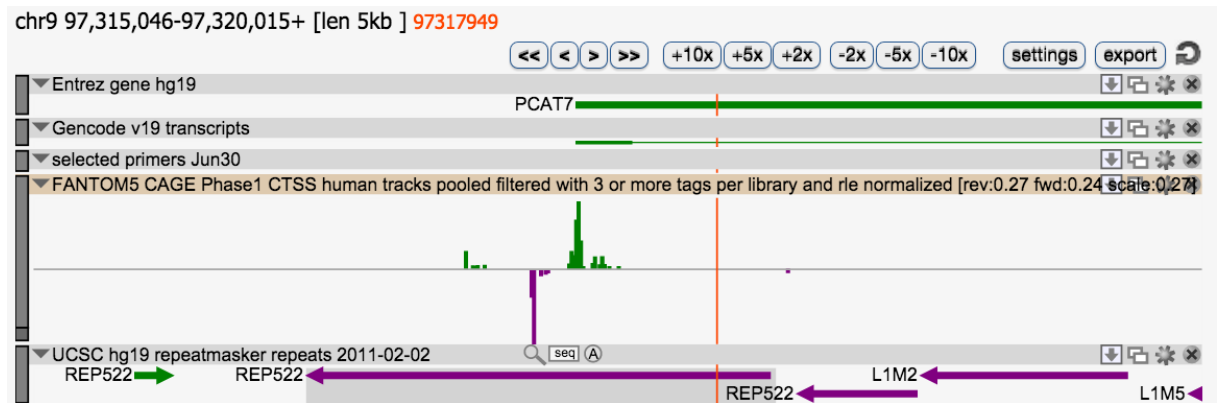


**BAGE2 - B melanoma antigen family, member 2**

1.8kb



**PCAT7 - prostate cancer associated transcript 7 (non-protein coding)**



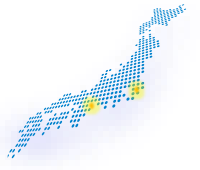
**REP522 without promoter**

**REP522 with promoter**

**REP522 without promoter**

<http://fantom.gsc.riken.jp/zenbu/>

# Summary



- Pan cancer biomarkers
  - protein coding
  - lncRNAs
  - enhancers
- Activation of transcription from repeats in cancer
- Specific activation of REP522 element

### Pan cancer paper

Yuji Tanaka

Hideya Kawaji

Albin Sandelin

Robin Andersson

Masayoshi Itoh

Timo Lassmann

Yoshihide Hayashizaki

Piero Carninci

Alistair Forrest

### LINC02021 KD

Jay Shin

Yuji Tanaka

Jasmine Ooi

Yuri Ishizu

Alistair Forrest

Piero Carninci

### FANTOM5 Consortium

#### Funding:

1. Japan Society for the Promotion of Science (JSPS Fellowship)
2. Foreign Postdoctoral Researcher program
3. Research Grant for RIKEN Omics Science Center from MEXT to YH.
4. Grant of the Innovative Cell Biology by Innovative Technology (Cell Innovation Program) from the MEXT, Japan to YH.
5. Research Grant from MEXT to the RIKEN Center for Life Science Technologies

Chung-Chau Hon

Michiel de Hoon

Jessica Severin

Giovanni Pascarella

Kosuke Hashimoto

Jordan Ramilowski

Erik Arner

**All DGT members for valuable discussions**



# Transcriptome Analysis of Recurrently Deregulated Genes Across Multiple Cancers Identifies New Pan-Cancer Biomarkers

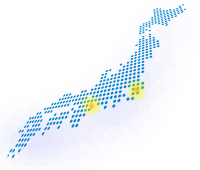
Bogumil Kaczkowski<sup>1</sup>, Yuji Tanaka<sup>1,2</sup>, Hideya Kawaji<sup>1,2,3</sup>, Albin Sandelin<sup>4</sup>, Robin Andersson<sup>4</sup>, Masayoshi Itoh<sup>1,3</sup>, Timo Lassmann<sup>1,5</sup>, the FANTOM5 consortium, Yoshihide Hayashizaki<sup>3</sup>, Piero Carninci<sup>1</sup>, and Alistair R.R. Forrest<sup>1,6</sup>

## Abstract

Genes that are commonly deregulated in cancer are clinically attractive as candidate pan-diagnostic markers and therapeutic targets. To globally identify such targets, we compared Cap Analysis of Gene Expression (CAGE) profiles from 225 different cancer cell lines and 339 corresponding primary cell samples to identify transcripts that are deregulated recurrently in a broad range of cancer types. Comparing RNA-seq data from 4,055 tumors and 563 normal tissues profiled in the The Cancer Genome Atlas and FANTOM5 datasets, we identified a core transcript set with theranostic potential. Our analyses also

revealed enhancer RNAs which are upregulated in cancer, defining promoters which overlap with repetitive elements (especially SINE/Alu and LTR/ERV1 elements) that are often upregulated in cancer. Lastly, we documented for the first time upregulation of multiple copies of the REP522 interspersed repeat in cancer. Overall, our genome-wide expression profiling approach identified a comprehensive set of candidate biomarkers with pan-cancer potential, and extended the perspective and pathogenic significance of repetitive elements which are frequently activated during cancer progression. *Cancer Res*; 1–11. ©2015 AACR.

Kaczkowski, B., et al. (2015). *Cancer Research*, ©2015 AACR



Thank you for your attention!

# REP522 -- largely palindromic, unclassified interspersed repeat



## Description

REP522 was originally called a telomeric satellite, even though in human, this 1.8Kb sequence is not found in tandem arrays, nor is it particularly restricted to telomeric regions. The central 1/3rd forms a >600bp imperfect palindrome. Other parts of these element display similarity to both LINE and LTR retrotransposons; to cut-and-paste DNA transposons; even to a Helitron! It is unclear how much of this is chimeric homology vs false positive.

The model is 1817 long.  
The average hit is 431.

