



ELSEVIER

Contents lists available at ScienceDirect

## Data in Brief

journal homepage: [www.elsevier.com/locate/dib](http://www.elsevier.com/locate/dib)

## Data Article

## Gene regulation by long purine tracks in brain related diseases

Himanshu Narayan Singh<sup>a,b</sup>, Moganty R. Rajeswari<sup>a,\*</sup><sup>a</sup> Department of Biochemistry, All India Institute of Medical Sciences, New Delhi 110029, India<sup>b</sup> School of Sciences, Noida International University, Gautam Budh Nagar, 203201 Uttar Pradesh, India

## ARTICLE INFO

## Article history:

Received 2 August 2015

Received in revised form

15 August 2015

Accepted 24 August 2015

Available online 4 September 2015

## Keywords:

Purine repeat

Human genome

Gene regulation

Brain disease

## ABSTRACT

Purine repeats are randomly distributed in the human genome, however, they show potential role in the transcriptional deregulation of genes. Presence of long tracks of purine repeats in the genome can disturb its integrity and interfere with the cellular behavior by introducing mutations and/or triple stranded structure formation in DNA. Our data revealed interesting finding that a majority of genes carrying purine repeats, of length  $n \geq 200$ , were down regulated and found to be linked with several brain related diseases [1]. The unique feature of the purine repeats found in the present study clearly manifests their significant application in developing therapeutics for neurological diseases.

© 2015 The Authors. Published by Elsevier Inc. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

## Specifications table

Subject area	Biology
More specific subject area	Genetics, Bioinformatics
Type of data	Table, Software generated files

DOI of original article: <http://dx.doi.org/10.1016/j.gene.2015.07.007>

\* Corresponding author.

E-mail address: [rajeswari3011@hotmail.com](mailto:rajeswari3011@hotmail.com) (M.R. Rajeswari).<http://dx.doi.org/10.1016/j.dib.2015.08.024>2352-3409/© 2015 The Authors. Published by Elsevier Inc. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

How data was acquired	Software generated
Data format	Analyzed
Experimental factors	Purine repeats ( $n \geq 200$ ) were searched in the human genome and also tried to explore their association with neurological disorders.
Experimental features	Purine repeat were searched by the help of home-made PERL script and further mapped them with neurological disorders
Data source location	New Delhi, India
Data accessibility	Data is supplied in this article

---

### Value of the data

---

- Identified purine repeats (PR,  $n \geq 200$ ) are unique in the human genome. Therefore, genes carrying purine repeats can be used as potential therapeutic tools in controlling gene expression and also in sequence-specific drug delivery.
  - The data will be helpful to explore the risk associated with acquiring disease causing mutations related to diseases.
  - The data will also be useful to study of evolutionary dynamics.
- 

## 1. Data, experimental design, materials and methods

### 1.1. Data resources

In present study, four data resources were utilized viz. (i) Human Genome Sequence: NCBI/Genome database; (ii) gene annotation: Ensemble Genome Browser; (iii) gene-disease association: GenAtlas database; and (iv) expression datasets: NCBI/GEO database. [Table 1](#)

## 2. Algorithm developed for purine repeat search

An indigenous PERL script “PuRepeatFinder.pl” was developed to locate PRs,  $n \geq 200$ , in the human genome. The tool enlists the PRs in the chronological order of its genomic coordinates along with PR-length and sequence. The script implements the knowledge based window-shift algorithm, and identify only uninterrupted, non-overlapping purine repeats.

## 3. Web tools

We have utilized two web-tools: (i) non-B DNA Motif Search Tool (nBMST): to search for the mirror repeat motifs with the identified PRs. It searches for the perfect and imperfect mirror repeats within the provided sequences [\[2\]](#) and (ii) Idiographica: to show the distribution of PR-genes on the chromosomes [\[3\]](#).

## 4. Microarray data analysis

Two open source R-packages of Bioconductor project viz. limma: used for agilent based microarray data, and affy: for affymatrix based microarray data, were used to calculate gene expression levels. Expression computation involves three steps: (i) background correction, (ii) normalization and (iii)

**Table 1**Description of PR-genes (polypurine nucleotides,  $n \geq 200$ ) associated with neurological disorders, PR sequences and its coordinates in human genome. PR: Purine repeat.

Gene symbol	Protein name	Contig	Chromosomal position		PR length	PR sequence
			Start	End		
RABGAP1L	RAB GTPase activating protein 1-like	NT_004487.19	26199819	26200018	200	AAAAAAAAAAGAAGAAGAAGAGGAA- GAGGAGGGGAGGGGGAGGAGGAGAAAGAAGAA- GAGGAGGAGGGGAGGGGGAGGAGGAAGAAAGAAGAGGAAGAGGA- GAGGAGGGGGAGGAGGAGAAAGAAGAA- GAGGAGGAGGGGGAGGGGGAGGAGGAGAAAGAAGAAGAA- GAAAAGGGGG
ALK	anaplastic lymphoma receptor tyrosine kinase	NT_022184.15	8875666	8876076	411	GAAGAAGAAGAAAAGAAGAAGAAGAAAAGAAGAAGAAAAGAAGAA- GAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAA- GAAGGGGAAGAGGAAGGGGAAGGAGGAGGAGGAGGAGAAGGAGAA- GAGGAAGGGGAGGAAGGGGGAGGAGGAGGAGGAGGAGGAGGAGGAG- GAGGAGAAGAAGAAGAAGAGGGGAAGAAGGGGAAGAAGGGGAA- GAAGGGGAAGAAGGAGAAGAGGAA- GAGGAAGGGGAAGGGGAAGGGGAAGGGGAAGGGGAAGGGAAGAGGAA- GAGGAAGAAGAAGAGGAAGAAGAAGGGAAGGGAAGGGAAGAAGAAGAA- GAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAA- GGGGAAGAAGGAAGGAAGGAAGGAAGGAGGAGGGGAA- GAAGGGGAAGGGGAAGGAGGGGAGGGGAAGGGAAGGGAAGGGAAG- GAAAGGAAAGGAAGGAA- GAAAGGGGAAGGGAAGGAAAGGAAAGGAAAGGGA- GAGGAAGGAAGGGAAAGGAAGGAAAGGGAAGGGAAGGGAAGGGAAGG- GAAGGAAGGGAAAGGAGGAAA- GAAAGGAAGGAAAGGAAGGAAAGGGAAGGGAAGGGAAGGGAAGGA- GAAGGAAGGGAAAGGGAAGGGA
GPR155	G protein-coupled receptor 155	NT_005403.17	25528765	25529048	284	AAAAAAAAAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGA- GAAAGAAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGA- GAAAGAGAAAGAGAAAAAGAAAGAAAGAGAAAGAAAGAAAAAGAGAAA- GAAAGAAGGAAAGAAAGAAAGAAAGAAAGAAAGAAAGAAAGAAAGAA- GAAAGAAAGAAAAAGAA
ROBO2	roundabout, axon guidance receptor, homolog 2 (Drosophila)	NT_022459.15	10576982	10577202	221	AAAGGAAGGAAGGAAGGAAAGAAAGAAAGAAAGAGAAAGGAGAGAGAAA- GAAAGAAAAGGGAAGGAAGGAAGGAAGGAAGGAAGGAAGGAAGAAAA- GAAAGAAAGAAAGAAAGAAAGGGAAGGAAGGAGAGAGAAAGAAAGGAAA- GAAAGAAAGAAAGAAAGAAAGAAAGAAAGAAAGAAAGAAAGAAAGAA- GAAAGAAAGAAAAAGAA
ARPP21	cAMP-regulated phosphoprotein, 21 kDa	NT_022517.18	35729714	35730024	311	AAAGGAAGGAAGGAAGGAAAGAAAGAAAGAAAGAAAGGAGAGAGAGAAA- GAAAGAAAGGGAAGGAAGGAAGGAAGGAAGGAAGGAAGGAAGGAAGAAA- GAAAGAAAGAAAGAAAGAAAGGGAAGGAAGGAGAGAGAAAGAAAGGAAA- GAAAGAAAGAAAGAAAGAAAGGGAAGGAAGGAGAGAGAAAGAAAGAA- GAAAGAAAGGAAAGAAAGAAAGAAAGAAAGAAAGAAAGAAAGGAGAGA- GAGGAGAGAGGGGAGGGAAGGAAGGAAGGAAGGAAGGAAGGAAGGAAG- GAGGAAAGAAAAG
			35649333	35649547	215	AAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAAGAA- GAGGAAGGGAAGAGGAAGGAAGGAGGAGGAGGAGGAGGAGGAGGAGAAGGA-









expression value computation [4]. Further, *t*-test was applied to screen statistically significant differential levels in mRNA expression of genes amongst patients and normal samples and  $p \leq 0.05$  were considered as significant [1].

## Acknowledgment

Himanshu Narayan Singh thanks the Indian Council of Medical Research, New Delhi, India for providing Senior Research Fellowship (3/1/2/43/Neuro/2013–NCD-1). Himanshu Narayan Singh is registered as Ph.D student at School of Sciences, Noida International University, Gautam Budh Nagar-203201, Uttar Pradesh, India.

## Appendix A. Supplementary material

Supplementary data associated with this article can be found in the online version at <http://dx.doi.org/10.1016/j.dib.2015.08.024>.

## References

- [1] H.N. Singh, M.R. Rajeswari, Role of long purine stretches in controlling the expression of genes associated with neurological disorders, *Gene* (2015), <http://dx.doi.org/10.1016/j.gene.2015.07.007>, pii: S0378-1119(15)00815-X.
- [2] R.Z. Cer, K.H. Bruce, D.E. Donohue, N.A. Temiz, U.S. Mudunuri, M. Yi, et al., Searching for non-B DNA-forming motifs using nBMST (non-B DNA motif search tool), *Current Protocols in Human Genetics*/Editorial Board Jonathan Haines L. Chapter 18 (2012) Unit 18.7.1–22. doi:10.1002/0471142905.hg1807s73.
- [3] T. Kin, Y. Ono, Idiographica: a general-purpose web application to build idiograms on-demand for human, mouse and rat, *Bioinformatics* 23 (2007) 2945–2946. <http://dx.doi.org/10.1093/bioinformatics/btm455>.
- [4] R.C. Gentleman, V.J. Carey, D.M. Bates, B. Bolstad, M. Dettling, S. Dudoit, et al., Bioconductor: open software development for computational biology and bioinformatics, *Genome Biol.* 5 (R80) (2004), <http://dx.doi.org/10.1186/gb-2004-5-10-r80>.