

Supplementary material for “Quantification of tumour evolution and heterogeneity via Bayesian epiallele detection”

James E. Barrett^{*1}, Andrew Feber¹, Javier Herrero¹, Miljana Tanic¹, Gareth Wilson^{2,3}, Charles Swanton^{2,3,4,†}, and Stephan Beck¹

[†]on behalf of the Lung TRACERx consortium

¹*UCL Cancer Institute, University College London, U.K.*

²*The Francis Crick Institute, London, U.K.*

³*Cancer Research U.K. Lung Cancer Centre of Excellence, UCL Cancer Institute, U.K.*

⁴*University College London Hospitals NHS Foundation Trust, U.K.*

June 2, 2017

A Extraction of viable loci

In Figure 1 is an example of an observed locus before any preprocessing steps have been taken. The locus is defined naturally since it does not overlap with any other observed sequencing reads. Note that several sequencing reads only partially overlap and many do not overlap at all. In addition, several CpG measurements are missing from the middle of some reads due to the paired-end sequencing protocol that was used to generate the data (sometimes the paired ends may not span the full length of the DNA fragment).

Intuitively this locus should be split into two loci as the two blocks of reads overlap by a single CpG which is not enough to phase the inferred epialleles. We implemented the following algorithm to split the observed sequencing reads into sensible loci.

1. Specify the minimum number of contiguous CpGs d_{min} and the minimum number of reads N_{min} required in order for a locus to be admissible (we used $d_{min} = 6$ and $N_{min} = 100$ in practice).
2. Specify the maximum proportion of missing data that is allowed (we used 25% in practice).
3. If the observed locus contains more than 25% missing data move to the next step, otherwise skip to step 4.
- 3(a). Represent the observed reads as a matrix with all non-missing measurements equal to 1 and all missing values represented with 0. Use hierarchical clustering to split the reads into two groups using the hamming distance.

^{*}Contact: regmjeb@ucl.ac.uk

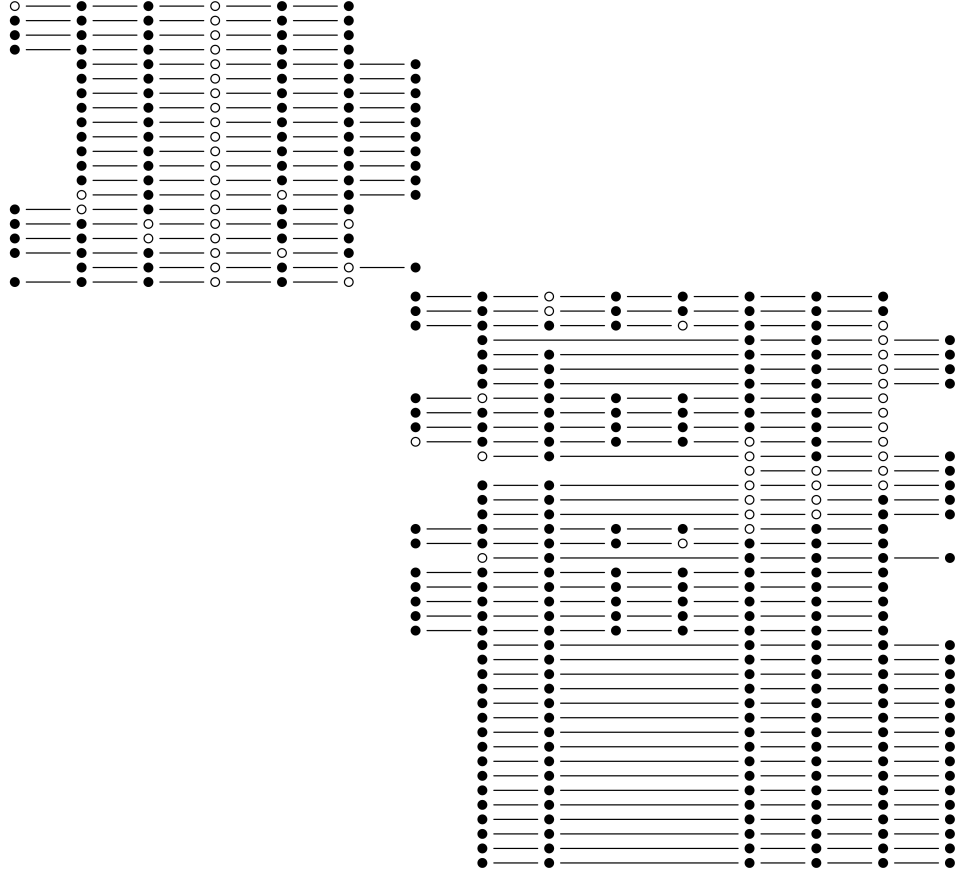


Figure 1: An example of an observed locus (chr1:15,232,224-15,232,587) before preprocessing.

- 3(b). If the two groups contain less than 25% of the missing data then proceed to the next step. Otherwise repeat the clustering above with three groups and so forth.
4. Discard any loci that fail to meet the minimum values of d and N .

We demanded that, in addition to the loci as a whole, each CpG site should not contain more than 25% missing data. Any CpGs that failed to meet this constraint were discarded. This helped to trim low-coverage CpGs from at the edges of observed loci (for example, the very leftmost CpG in Figure 1). The resulting loci after preprocessing are shown in Figure 2.

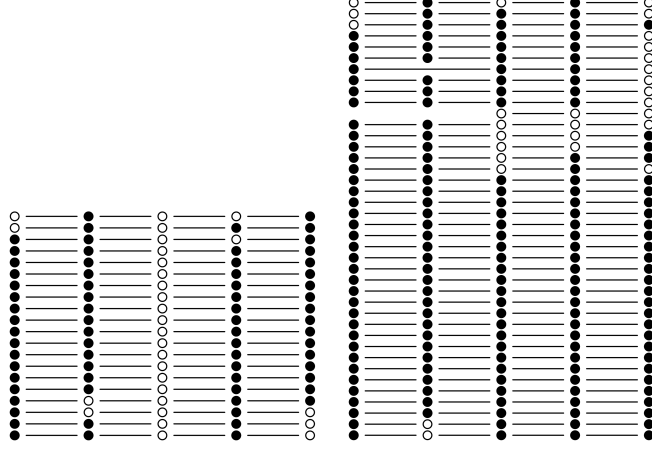


Figure 2: The two loci that are extracted from the observed data in Figure 1.

B MAP estimate for ϵ

Given \mathbf{X} , \mathbf{w} , Q and uniform priors $p(\mathbf{X}|Q)$ and $p(\mathbf{w}|Q)$ the MAP estimate for the hyperparameter ϵ is given by maximising the log of the posterior (2):

$$\begin{aligned}
 \mathcal{L}(\epsilon) &= \log p(\mathbf{Y}|\mathbf{X}, \mathbf{w}, \epsilon, Q) \\
 &= \log \prod_{i=1}^N \prod_{\mu=1}^d p(y_{i\mu}|x_{w_i\mu}, \epsilon, Q) \\
 &= \alpha_0 \log \epsilon + \alpha_1 \log(1 - \epsilon).
 \end{aligned} \tag{1}$$

Note that on the second line that if read \mathbf{y}_i originates from epiallele q then $w_i = q$. Recall that $\alpha_1 = \sum_{i,\mu} \delta_{y_{i\mu}, x_{w_i\mu}}$ and $\alpha_0 = \sum_{i,\mu} 1 - \delta_{y_{i\mu}, x_{w_i\mu}}$ denote the total number of matches and mismatches between the observed reads \mathbf{y} and the corresponding epialleles \mathbf{x} at this particular loci. Solving $d\mathcal{L}/d\epsilon = 0$ yields $\epsilon = \alpha_0/(\alpha_0 + \alpha_1)$.

C Simulation results

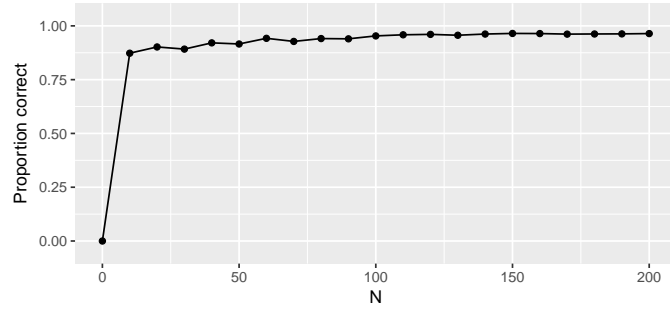


Figure 3: Proportion of observed reads attributed to the correct underlying epiallele as a function of N (the number of sequencing reads at the simulated locus). Parameters were fixed to $\epsilon = 0.05$, $d = 6$ and $Q = 3$.

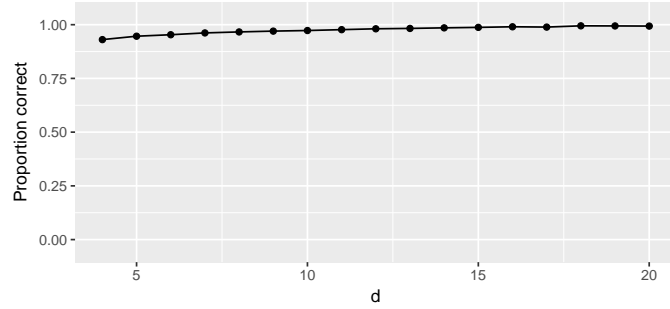


Figure 4: Proportion of observed reads attributed to the correct underlying epiallele as a function of d (the number of CpGs at the simulated locus). Parameters were fixed to $N = 100$, $\epsilon = 0.05$ and $Q = 3$.

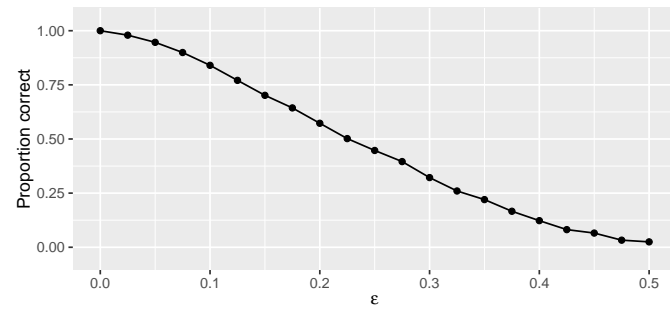


Figure 5: Proportion of observed reads attributed to the correct underlying epiallele as a function of ϵ (the noise level). Parameters were fixed to $N = 100$, $d = 6$ and $Q = 3$.

D Purity estimation

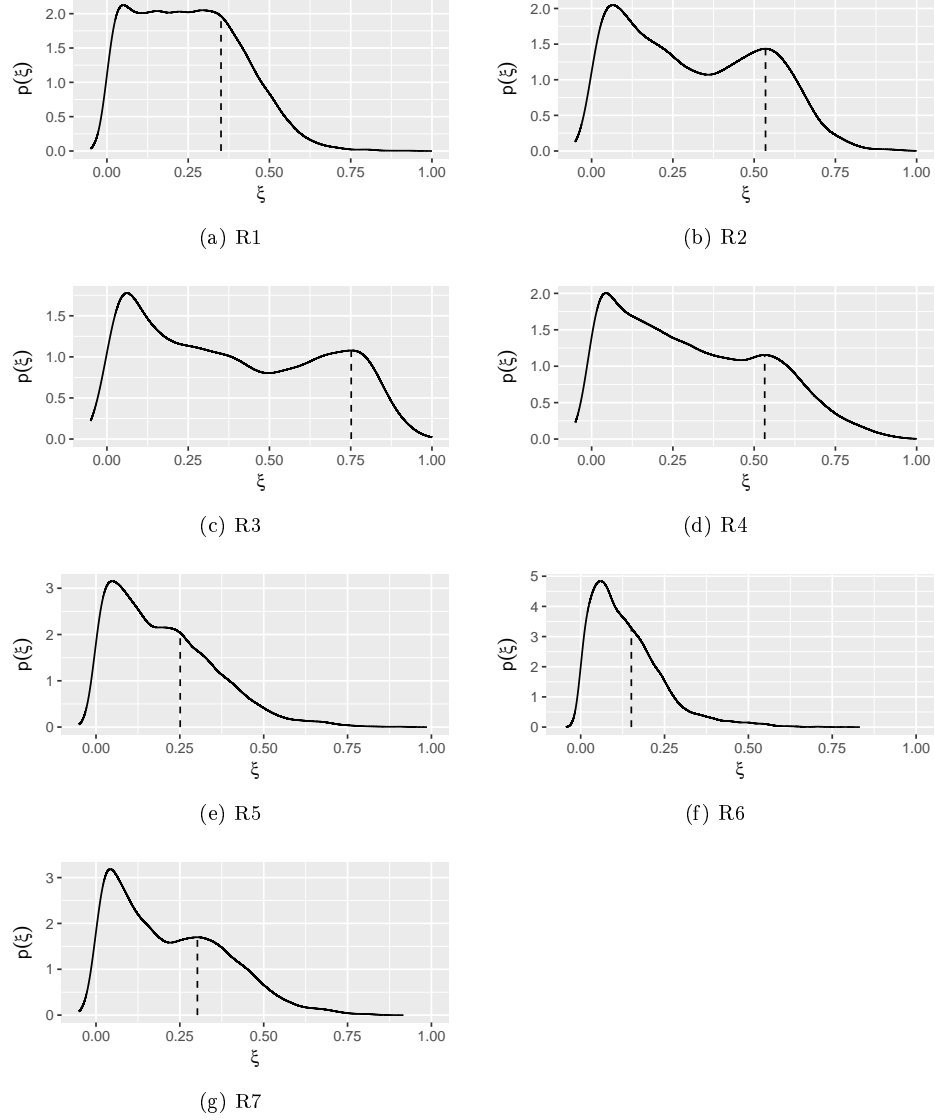


Figure 6: Empirical density plots of ξ , the proportion of epialleles at a locus that are different from normal tissue. The distribution of ξ will depend on the tumour purity since samples that are contaminated with less normal tissue will exhibit a greater deviance from the matched normal tissue epialleles. On this basis, the rightmost maxima (marked with a dashed vertical line) of the empirical densities are interpreted as a proxy for sample purity.

E Supplementary figures

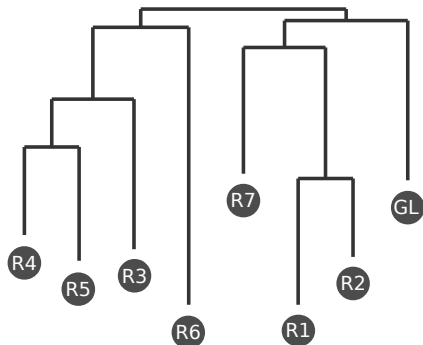


Figure 7: Phylogenetic tree generated from exome sequencing data from the same tumour that is studied in the main text. The exome data were generated and analysed independently as part of the Jamal-Hanjani et al., 2017 study. GL denotes germline.

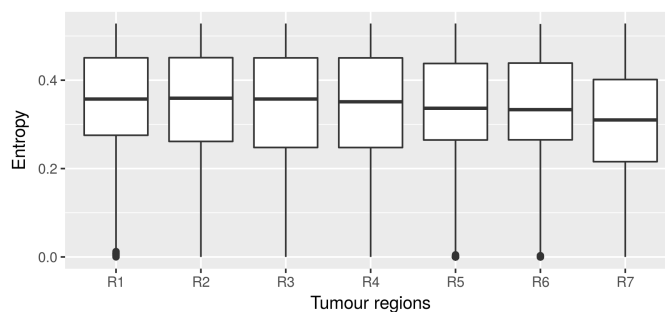


Figure 8: Box plots of the Shannon entropy of the epiallele distribution across the seven tumour regions (R1-R7) after decontamination of normal tissue.

F Comparison to alternative measures of disorder

As discussed in the main text the *epiallele entropy* provides a measure of epigenetic disorder within each tissue sample. In this section we compare our epiallele entropy approach to previously proposed measures of epigenetic disorder. In particular, we will compare the *epipolymorphism score* from Landan et al. (2012), the *methylation entropy* from Xie et al. (2011) and the *proportion of discordant reads* (PDR) approach from Landau et al. (2014).

Both the epipolymorphism score and the PDR score are restricted to loci with 4 CpGs. We have therefore restricted this comparison to only loci with 4 CpGs of which there are 38,831 in total across all of our tissue samples. Furthermore, none of the comparators can handle missing data so any reads with a missing CpG measurement were discarded. As described in the main text for our

epiallele entropy approach we used the distribution of epialleles after marginalisation over the \mathbf{w} parameter and any epialleles that had a frequency of less than 5% were discarded.

Box plots comparing the four different measures are shown in Figure 9. All four methods attribute lower disorder scores to the normal tissue. Low purity tumour regions such as R5, R6 and R7 are also given lower disorder scores by all four approaches. The greatest difference between normal and tumour tissue is observed using our epiallele entropy measure. Our method has the additional advantage that it can handle an arbitrary number of CpGs per locus and missing data.

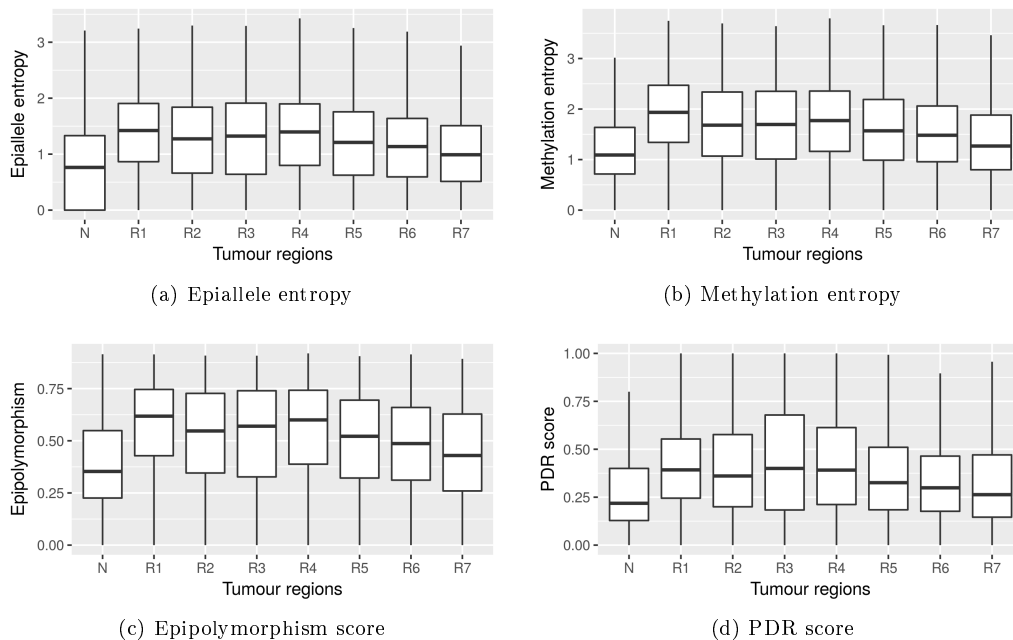


Figure 9: Box plots comparing different epigenetic disorder measures.

G Experimental datasets

The RRBS data are available at the European Nucleotide Archive under accession numbers ERS1546024, ERS1546025 and ERS1546026.

References

- G. Landan et al. Epigenetic polymorphism and the stochastic formation of differentially methylated regions in normal and cancerous tissues. *Nat. Genet.*, 44(11):1207–1214, 2012.
- D. Landau et al. Locally disordered methylation forms the basis of intratumor methylome variation in chronic lymphocytic leukemia. *Cancer Cell*, 26(6):813–825, 2014.

H. Xie et al. Genome-wide quantitative assessment of variation in DNA methylation patterns.
Nucleic Acids Res., 39(10):4099–4108, 2011.

H TRACERx consortium members

The TRACERx study (Clinicaltrials.gov no: NCT01888601) is sponsored by University College London (UCL/12/0279) and has been approved by an independent Research Ethics Committee (13/LO/1546). TRACER is funded by Cancer Research UK (grant number C11496/A17786) and coordinated through the Cancer Research UK & UCL Cancer Trials Centre.

Consortium members

Charles Swanton^{1,2,5}, Mariam Jamal-Hanjani¹, Selvaraju Veeriah¹, Seema Shafi¹, Justyna Czyzewska-Khan¹, Diana Johnson¹, Joanne Laycock¹, Leticia Bosshard-Carter¹, Gerald Goh¹, Rachel Rosenthal¹, Pat Gorman¹, Nirupa Murugaesu¹, Robert E Hynds^{1,3}, Gareth Wilson^{1,2}, Nicolai J Birkbak^{1,2}, Thomas B K Watkins², Nicholas McGranahan^{1,2}, Stuart Horswell², Richard Mitter², Mickael Escudero², Aengus Stewart², Peter Van Loo², Andrew Rowan², Hang Xu², Samra Turajlic^{2,4}, Crispin Hiley², Christopher Abbosh¹, Jacki Goldman², Richard Kevin Stone², Tamara Denner², Nik Matthews², Greg Elgar², Sophia Ward², Jennifer Biggs², Marta Costa², Sharmin Begum², Ben Phillimore², Tim Chambers², Emma Nye², Sofia Graca², Maise Al Bakir², Kroopa Joshi¹, Andrew Furness¹, Assma Ben Aissa¹, Yien Ning Sophia Wong¹, Andy Georgiou¹, Sergio Quezada¹, John A Hartley¹, Helen L Lowe¹, Javier Herrero¹, David Lawrence⁵, Martin Hayward⁵, Nikolaos Panagiotopoulos⁵, Shyam Kolvekar⁵, Mary Falzon⁵, Elaine Borg⁵, Teresa Marafioti⁵, Celia Simeon⁵, Gemma Hector⁵, Amy Smith⁵, Marie Aranda⁵, Marco Novelli⁵, Dahmane Oukrif⁵, Sam M Janes⁵, Ricky Thakrar⁵, Martin Forster⁵, Tanya Ahmad⁵, Siow Ming Lee⁵, Dionysis Papadatos-Pastos⁵, Dawn Carnell⁵, Ruheena Mendes⁵, Jeremy George⁵, Neal Navani⁵, Asia Ahmed⁵, Magali Taylor⁵, Junaid Choudhary⁵, Yvonne Summers⁶, Raffaele Califano⁶, Paul Taylor⁶, Rajesh Shah⁶, Piotr Krysiak⁶, Kendadai Rammohan⁶, Eustace Fontaine⁶, Richard Booton⁶, Matthew Evison⁶, Phil Crosbie⁶, Stuart Moss⁶, Faiza Idries⁶, Leena Joseph⁶, Paul Bishop⁶, Anshuman Chaturved⁶, Anne Marie Quinn⁶, Helen Doran⁶, Angela leek⁷, Phil Harrison⁷, Katrina Moore⁷, Rachael Waddington⁷, Juliette Novasio⁷, Fiona Blackhall⁸, Jane Rogan⁷, Elaine Smith⁶, Caroline Dive⁹, Jonathan Tugwood⁹, Ged Brady⁹, Dominic G Rothwell⁹, Francesca Chemi⁹, Jackie Pierce⁹, Sakshi Gulati⁹, Babu Naidu¹⁰, Gerald Langman¹⁰, Simon Trotter¹⁰, Mary Bellamy¹⁰, Hollie Bancroft¹⁰, Amy Kerr¹⁰, Salma Kadiri¹⁰, Joanne Webb¹⁰, Gary Middleton¹⁰, Madava Djearaman¹⁰, Dean Fennell¹¹, Jacqui A Shaw¹¹, John Le Quesne¹¹, David Moore¹¹, Apostolos Nakas¹², Sridhar Rathinam¹², William Monteiro¹³, Hilary Marshall¹³, Louise Nelson¹², Jonathan Bennett¹², Joan Riley¹², Lindsay Primrose¹², Luke Martinson¹², Girija Anand¹⁴, Sajid Khan¹⁵, Anita Amadi¹⁶, Marianne Nicolson¹⁷, Keith Kerr¹⁷, Shirley Palmer¹⁷, Hardy Remmen¹⁷, Joy Miller¹⁷, Keith Buchan¹⁷, Mahendran Chetty¹⁷, Lesley Gomersall¹⁷, Jason Lester¹⁸, Alison Edwards¹⁸, Fiona Morgan¹⁹, Haydn Adams¹⁹, Helen Davies¹⁹, Malgorzata Kornaszewska²⁰, Richard Attanoos²¹, Sara Lock²², Azmina Verjee²², Mairead MacKenzie²³, Maggie Wilcox²³, Harriet Bell²⁴, Natasha Iles²⁴, Allan Hackshaw²⁴, Yenting Ngai²⁴, Sean Smith²⁴, Nicole Gower²⁴, Christian Ottensmeier²⁵, Serena Chee²⁵, Benjamin Johnson²⁵, Aiman Alzetani²⁵, Emily Shaw²⁵, Eric Lim²⁶, Paulo De Sousa²⁶, Monica Tavares Barbosa²⁶, Alex Bowman²⁶, Simon Jordan²⁶, Alexandra Rice²⁶, Hilgardt Raubenheimer²⁶, Chiara Proli²⁶, Maria Elena Cufari²⁶, John Carlo Ronquillo²⁶, Angela Kwayie²⁶, Harshil Bhayani²⁶, Morag Hamilton²⁶, Yusura Bakar²⁶, Natalie Mensah²⁶, Lyn Ambrose²⁶, Anand Devaraj²⁶, Silviu Buder²⁶, Jonathan Finch²⁶, Leire Azcarate²⁶, Hema Chavan²⁶, Sophie Green²⁶, Hillaria Mashinga²⁶, Andrew G Nicholson^{26, 27}, Kelvin Lau²⁸, Michael Sheaff²⁸, Peter Schmid²⁸, John Conibear²⁸, Veni Ezhil²⁹, Babikir Ismail²⁹, Melanie Irvin-sellers²⁹, Vineet Prakash²⁹, Peter

Russell³⁰, Teresa Light³⁰, Tracey Horey³⁰, Sarah Danson³¹, Jonathan Bury³¹, John Edwards³¹, Jennifer Hill³¹, Sue Matthews³¹, Yota Kitsanta³¹, Kim Suvarna³¹, Patricia Fisher³¹, Allah Dino Keerio³¹, Michael Shackcloth³², John Gosney³², Pieter Postmus³², Sarah Feeney³², Julius Asante-Siaw³², Tudor Constatin³³, Raheleh Salari³³, Nicole Sponer³³, Ashwini Naik³³, Bernhard Zimmermann³³, Hugo J.W.L. Aerts³⁴, Stefan Dentre³⁵, Christophe Dessimoz^{36,37,38}.

Affiliations

1. Cancer Research UK Lung Cancer Centre of Excellence, University College London Cancer Institute, United Kingdom
2. The Francis Crick Institute, United Kingdom
3. Lungs for Living, UCL Respiratory, University College London, United Kingdom
4. The Royal Marsden Hospital, United Kingdom
5. University College London Hospitals NHS Foundation Trust, United Kingdom
6. University Hospital of South Manchester, United Kingdom
7. Manchester Cancer Research Centre Biobank, United Kingdom
8. Christie NHS Foundation Trust, Manchester, United Kingdom
9. Cancer Research UK Manchester Institute, United Kingdom
10. Heart of England NHS Foundation Trust, Birmingham, United Kingdom
11. Cancer Studies and Molecular Medicine, University of Leicester, United Kingdom
12. Leicester University Hospitals, United Kingdom
13. National Institute for Health Research Leicester Respiratory Biomedical, Research Unit, United Kingdom
14. North Middlesex Hospital, United Kingdom
15. Royal Free Hospital, United Kingdom
16. Barnet Hospital, United Kingdom
17. Aberdeen Royal Infirmary, United Kingdom
18. Velindre Cancer Centre, Cardiff, Wales, United Kingdom
19. Cardiff & Vale University Health Board, Cardiff, Wales, United Kingdom
20. University Hospital Of Wales Heath Park, Cardiff, Wales, United Kingdom
21. Department of Pathology, University Hospital of Wales and Cardiff University, Heath Park, Cardiff, Wales, United Kingdom
22. The Whittington Hospital NHS Trust, United Kingdom
23. Independent Cancer Patients Voice, United Kingdom
24. Cancer Research UK & UCL Cancer Trials Centre, United Kingdom
25. University Hospital Southampton NHS Foundation Trust, United Kingdom
26. Royal Brompton and Harefield NHS Foundation Trust, United Kingdom
27. National Heart and Lung Institute, Imperial College, United Kingdom
28. Barts Health NHS Trust, United Kingdom
29. Ashford and St. Peter's Hospitals NHS Foundation Trust, United Kingdom
30. The Princess Alexandra Hospital NHS Trust, United Kingdom
31. Sheffield Teaching Hospitals NHS Foundation Trust, United Kingdom
32. Liverpool Heart and Chest Hospital NHS Foundation Trust, United Kingdom
33. Natera Inc., 201 Industrial Road, Suite 410, San Carlos, CA 94070
34. Dana-Farber Cancer Institute, Brigham & Women's Hospital, Harvard Medical School, 450 Brookline Ave, JF518, Boston, MA 02115-5450, USA

35. Wellcome Trust Sanger Institute, Hinxton, CB10 1SA, United Kingdom
36. Bioinformatics Group, Department of Computer Science, University College London
37. University of Lausanne
38. Swiss Institute of Bioinformatics