1    Universal latent axes capturing Parkinson's patient deep phenotypic variation reveals patients

2                    with a high genetic risk for Alzheimer's disease are more likely

3                          to develop a more aggressive form of Parkinson's.

4    Cynthia Sandor DMV,PhD[1†], Stephanie Millin DPhil[2†], Andrew Dahl DPhil[3], Michael

5    Lawton PhD[4], Leon Hubbard PhD[5], Bobby Bojovic MS[3], Marine Peyret-Guzzon PhD[3],

6    Hannah Matten MS[3], Christine Blancher PhD[3], Nigel Williams PhD[5], Yoav Ben-Shlomo

7    PhD[4], Michele T. Hu MD, PhD[6,7], Donald G. Grosset MD, PhD[8], Jonathan Marchini PhD[3,9,*],

8    Caleb Webber PhD[1,2,*]

9

10   **Affiliations:**

11   [1]UK Dementia Research Institute, Cardiff University, Cardiff, UK

12   [2]Department of Physiology, Anatomy and Genetics, University of Oxford, Oxford, UK

13   [3]Wellcome Centre for Human Genetics, University of Oxford, Oxford, UK

14   [4]School of Social and Community Medicine, University of Bristol, Bristol, UK

15   [5]MRC Centre for Neuropsychiatric Genetics and Genomics, Institute of Psychological

16   Medicine and Clinical Neurosciences, School of Medicine, Cardiff University, Cardiff, UK

17   [6]Oxford Parkinson's Disease Centre, Department of Physiology, Anatomy and Genetics, Le

18   Gros Clark Building, University of Oxford, Oxford, UK

19   [7]Nuffield Department of Clinical Neurosciences, Division of Clinical Neurology, University

20   of Oxford, Oxford, UK.

21   [8]Department of Neurology, Institute of Neurological Sciences, Queen Elizabeth University

22   Hospital, Glasgow, United Kingdom

1    [9]Department of Statistics, University of Oxford, Oxford, UK

2

3    *To whom correspondence should be addressed:  E-mail:  webberc4@cardiff.ac.uk &

4    marchini@stats.ox.ac.uk

5    †Both authors contributed equally to this work.

6

# 1  Abstract

2   The generation of deeply phenotyped patient cohorts offers an enormous potential to identify

3   disease subtypes but are currently limited by the cohort size and the heterogeneity of the

4   clinical assessments collected across different cohorts. Identifying the universal axes of clinal

5   severity and progression is key to accelerating our understanding of how disease manifests

6   and progresses. These universal axes would accelerate our understanding of how Parkinson's

7   disease (PD) manifests and progresses through which patients may be appropriately

8   compared appropriately stratified, and personalised therapeutic strategies and treatments can

9   be developed and targeted. We developed a Bayesian multiple phenotype mixed model

10   incorporating the genetic relationships between individuals which is able to reduce a wide-

11   array of different clinical measurements into a smaller number of continuous underlying

12   factors named phenotypic axis. We identify three principal axes of PD patient phenotypic

13   variation which are reproducibly found across three independent, deeply and diversely

14   phenotyped cohorts. Together they explain over 75% of the observed clinical variation and

15   remain robustly captured with a fraction of the clinically-recorded features. The most

16   influential axis was associated with the genetic risk of Alzheimer's disease (AD) and involves

17   genetic pathways associated with neuroinflammation. Our results suggest PD patients with a

18   high genetic risk for AD are more likely to develop a more aggressive form of PD including,

19   but not limited to, dementia.

20

## 1   Introduction

2      A critical challenge in medicine is to understand why the clinical presentations of

3   each patient affected by the same disorder vary. This is especially true for Parkinson's disease

4   (PD), for which the age of onset, the rate of progression, type and severity of symptoms differ

5   across more than a million people worldwide living with this disease [1]. To accelerate the

6   identification of disease subtypes, large deeply phenotyped cohorts of PD patients have been

7   created, in which valuable clinical, imaging, biosample and genetic data has been collected,

8   and increasingly with longitudinal monitoring [2-4].

9      Recent studies exploiting these deeply phenotyped cohorts have classified patients

10   into discrete phenotypic subgroups, each displaying a characteristic set of symptoms [5-7]. To

11   define PD subtypes, most of these studies employ some form of variable selection to create a

12   distance matrix between individuals, followed by clustering methods such as k-means or

13   hierarchical clustering. These methods provide discrete phenotypic groups, which are

14   appealing in their categorical nature but have many shortfalls. Firstly, while selection

15   methods quantify how much variance each phenotype explains, no robust method was used to

16   define a threshold for this measure above which a phenotype contributes to the distance

17   matrix. Consequently, the definition of which phenotypes are essential to group patients and

18   which are irrelevant can be somewhat arbitrary. For example, two recent studies [5, 8], using the

19   same Parkinson's Progression Markers Initiative (PPMI) cohort show divergent results:

20   apathy and hallucinations were key subtype classifiers in the first study [8], but not in the

21   second one [5], because these variables were not included. Secondly, K-means clustering

22   requires the number of phenotypic groups to be prespecified, and this choice has the potential

23   to be biased towards preconceived expectations with smaller groups ignored or erroneously

24   joined with larger groups.   Finally, the creation of discrete groups may not reflect the

1    possibly continuous nature of phenotypic variability and ignores the greater statistical power

2    of continuous traits.

3          To overcome these limitations, we propose here an approach focused on the

4    continuous variation of phenotypes. Rather than focusing on presence versus absence, or mild

5    versus severe phenotypes, we incorporate the whole spectrum of severity displayed across the

6    population. For this, we applied PHENIX (PHENotype Imputation eXpediated), a multiple

7    phenotype mixed model (MPMM) approach initially developed to impute missing

8    phenotypes [9], that can also be exploited for genetically-guided dimensionality reduction of

9    multiple traits. This approach models the phenotypes as a combination of genetic and

10   environmental factors and the genetic component is computed from the correlation matrix

11   between the individual's genetic data.

12         Applying PHENIX to the deeply phenotyped UK-based *Discovery* cohort, we identify

13   a small number of axes underlying individual PD patient phenotypic variation that explain the

14   variation in the much larger number of clinically-observed phenotypes. We demonstrate the

15   universality of these axes of phenotypic variation amongst PD patients by independently

16   deriving similar axes in each of three cohorts: UK *Tracking* cohort including 1807

17   individuals, the UK *Discovery* cohort including 842 PD patients and US PPMI cohort

18   including 439 PD patients that has a different clinical structure from the UK cohorts. We

19   show that this reproducibility is not achieved by other commonly-used dimensionality-

20   reduction methods. Finally, we demonstrate that the most influential axis was associated with

21   the genetic risk of Alzheimer's disease (AD) suggesting PD patients with a high genetic risk

22   for Alzheimer's disease are more likely to develop a more aggressive form of PD including

23   dementia symptoms.

24

## 1    Materials and Methods

### 2    Discovery cohort

3    We considered 842 PD cases from the *Discovery* cohort constituted of 1700 subjects,

4    including over 1000 people with Parkinson's, plus 320 healthy controls and 340 individuals

5    thought to be 'at-risk' of developing future Parkinson's. Individuals were required to have at

6    least 90% chance of PD according to UK-Parkinson's disease brain bank criteria, no

7    alternative diagnosis and disease duration less than 3.5 years. All patients have a clinical

8    assessment repeated every eighteen months and have been already described[4,6]. Phenotype

9    data were collected for over a hundred clinical attributes, affecting autonomic, neurological

10   and motor phenotypes (**Supplementary Fig. 1**) and described in the **Supplementary Table**

11   **1.** Genotype data were generated using the Illumina HumanCoreExome-12 v1.1 and Illumina

12   InfiniumCoreExome-24 v1.1 SNP arrays.

### 13   UK Tracking Parkinson's study

14        We considered 1807 PD cases from the *Tracking* Parkinson's cohort, which was

15   already described in detail by Malek *et al.* [2]  and was used to identify the impact of mutations

16   within glucocerebrosidase gene (*GBA*) on different PD clinicals manifestations [10]. Genotype

17   data were generated using the Illumina Human Core Exome array.

### 18   PPMI cohort

19        The PPMI cohort (http://www.ppmi-info.org) was already described in detail

20   (including PPMI protocol of recruitment and informed consent) by Marrek *et al.* [11]. We

21   downloaded data from the PPMI database on September 2017 in compliance with the PPMI

22   Data Use Agreement. We considered 472 newly-diagnosed typical PD subjects: subjects with

23   a diagnosis of PD for two years or less and who are not taking PD medications. We used the

24   baseline (t=0) of clinical assessments, described in detail in the **Supplementary Table 2**. We

1    excluded any individual with > 5% of missing data (437 individuals included). Participants

2    have been genotyped using two genotyping arrays, ImmunoChip [12] and NeuroX [13], [14]. As

3    more participants were genotyped on NeuroX array, we used the genotype data of the

4    NeuroX chip.

5    **Methods**

6    **Genotype: quality control & Imputation**

7    Quality control was carried out independently using PLINK v1.9 [15] (SI). Imputation of

8    unobserved and missing variants was carried out separately for each cohort (SI)

9    **Phenotypic axis**

10    Our continuous measures of severity are based on a multiple phenotypes mixed model

11    approach (MPMM) named PHENIX (PHENotype Imputation eXpediated) which includes

12    genetic relationships between individuals, and is designed to impute missing phenotypes [9].

13    To impute missing phenotypes, PHENIX reduces the variation within a cohort to a smaller

14    number of underlying factors that are then used to predict individual missing values. Here,

15    we exploit the identification of these underlying factors as providing the latent axes of patient

16    variation which underlie a larger number of clinically observed phenotypes. The outcome is

17    that the many clinical phenotypes (sometimes missing for some individuals) of each

18    individual are represented through a smaller number of underlying latent variables of

19    phenotypic variation that manifest the observed clinical phenotypes, which we name herein

20    as *phenotypic axes*.

21    PHENIX [9] use a Bayesian multiple-phenotype mixed model (MPMM), where the correlations

22    between clinical phenotypes (Y) are decomposed into a genetic and a residual component

23    with the following model: $Y=U+e$, where U represents the aggregate genetic contribution

24    (whole genotype) to phenotypic variance and e is idiosyncratic noise. As the estimation of

1    maximum likelihood covariance estimates can become computationally expensive with

2    increasing number of phenotypes, PHENIX uses a Bayesian low-rank matrix factorization

3    model for the genetic term U such as: $U = S\beta$, in which $\beta$ is can be used to estimate the

4    genetic covariance matrix between phenotypes and S represents a matrix of latent

5    components that each follow $\sim N(0,G)$ where G is the Estimate of Relatedness Matrix from

6    genotypes. The resulting latent traits (S) are used as phenotypic axes, each representing the

7    severity of a number of non-independent clinical phenotypes. The details to run PHENIX and

8    extract the phenotypic axes are given in the **Supplemental Information.**

9    **Disease phenotypic axis**

10   We derived disease phenotypic axis consisting to replace the general genetic component in a

11   MPMM by a disease risk genetics component. To calculate a disease relatedness matrix, we

12   considered only genetic variants (after pruning) associated with human complex traits. For

13   different complex human traits with GWA results publically available (**Supplementary**

14   **Table 01**), we calculated a disease relatedness, that we used subsequently to derive

15   phenotypic axes (**Supplementary Information**).

16

17   **Results**

18   **Three continuous measures capture 75% of the clinical variation.**

19   Initially, we generated phenotypic axes from a cohort of 842 PD patients (*Discovery* cohort)

20   which had been genotyped and phenotypically characterised with 40 clinical assessments

21   (**Supplementary Table 1**). Each latent axis reflected a number of co-varying observed

22   clinical assessments. Among the phenotypic axes that explained more than 5%, Axes 1, 2 and

23   3 explained 39.6%, 28.7% and 6.8% of the clinical variation respectively. Together, these 3

24   top axes account for over 75% of the clinically-observed variation (**Supplementary Fig. 2**).

1    To examine whether similar phenotypic axes are obtained in different deeply phenotyped PD

2    cohorts, we derived phenotypic axes within an independent cohort of 1807 PD individuals

3    from the UK *Tracking* cohort [2] that had made similar clinical observations to the *Discovery*

4    cohort. We found significant Pearson's correlation coefficients between each cohort's first

5    three phenotypic axes: Axis 1 r=0.92 (p=3 x $10^{-13}$), Axis 2 r=0.89 (p=4 x $10^{-11}$), Axis 3

6    r=0.72 (p=5 x $10^{-6}$) (**Fig. 1**). Nevertheless, a major concern was that the identification of the

7    same phenotypic axes might, at least in part, be due to the very similar structure of the

8    clinical phenotyping between the two UK cohorts. To address this, we examined the

9    independent US-based PPMI cohort consisting of 439 sporadic PD individuals that had been

10   clinically phenotyped following a substantially different protocol to the UK cohorts. After

11   deriving phenotypic axes in the PPMI cohort, we found significant similarities between the

12   first three phenotypic axes derived for both the *Discovery*-UK and PPMI-US cohorts: the

13   coefficients of determination (R^2) between three first axes across different categories of

14   clinical phenotypes from each cohort were: Axis1: 0.665 (p=0.048), Axis 2: 0.914 (p=0.003)

15   and Axis 3: 0.754 (p=0.025) (**Fig. 2 & Supplementary Figure 3**). These consistent

16   similarities in the axes of phenotypic variation independently derived for each of three

17   different PD cohorts demonstrates the reproducibility of these axes of phenotypic variation

18   amongst Parkinson's patients. Finally, by comparing PHENIX with other methods of

19   dimensionality    reduction,    specifically    Principle    Component    Analyses    (PCA),

20   Multidimensional Scaling (MDS) and Independent component analysis (ICA), only the

21   dimensions discovered by the MPMM model, PHENIX, were significantly correlated

22   between both cohorts and thus no other method was able to identify similar axes of

23   phenotypic variation across UK and US PD cohorts (**Fig. 2**).

24   **Each phenotypic axis represents a distinct set of clinical features**

1       To interpret the clinical relevance of each phenotypic axis, we examined the

2   correlation between individual clinical features and the phenotypic axes (**Table 1 & Fig. 1 &**

3   **Supplementary Figure 4**). We observed that each phenotypic axis corresponded to a subset

4   of clinical features, differing in both extents and directions of severity. Axis 1 represented

5   worsening non-tremor motor phenotypes, anxiety and depression accompanied by a decline

6   of the cognitive function (**Table 1 & Fig. 3**). Worsening anxiety and depression were also

7   features of Axis 2, in addition to increasing severity of autonomic symptoms and increasing

8   motor dysfunction. Axis 3 was associated with general motor symptom severity including

9   rigidity, bradykinesia and tremor of the whole body independently of non-motor features.

10   The contribution of different phenotypes to these axes was therefore highly variable. Specific

11   aspects of motor dysfunction were important factors in defining the majority of axes. Anxiety

12   and depression were also relatively important features, but only for axes explaining the

13   largest amounts of variation. Conversely, cognitive impairment was associated only with

14   Axis one. However, this observation must be weighted by the fact that cognitive

15   impairment/dementia are reported at a later disease stage and thus likely under-represented in

16   recently diagnosed cases.

17       Although each phenotypic axis is associated with a distinct set of clinical features,

18   they are not independent but instead strongly correlated (**Supplementary Figure 5**). We find

19   no significant relation between the phenotypic axes and principal components of genetic

20   ancestry (**Methods**) suggesting that the phenotypic axes are not biased by the population

21   structure (**Supplementary Figure 5, Supplementary Table 3**). However, as previously

22   reported, gender influences clinical symptoms [4] and we also observe a significant association

23   between gender and Axis 2 (**Supplementary Table 3,** p=4.5x10$^{-5}$).

24   To assess to what extent the phenotypic axes might be affected by the number of clinical

25   observations, within the *Discovery* cohort we compared the phenotypic axes built on all

1  clinical features with phenotypic axes generated with incomplete sets of randomly-selected

2  clinical features. We observed a strong correlation (r > 0.8) between each of the two first

3  phenotypic axes built with as few as 50% of the clinical variables and their respective

4  original phenotypic axes, suggesting that these two axes are extremely robust in terms of the

5  numbers of clinical variables considered (**Supplementary Figure 6**).

6  **The integration of genetic relationships between patients improves capture of the**

7  **Parkinson's disease clinical variation and reproducibility.**

8         The PHENIX MPMM approach employed here to derive phenotypic axes exploits the

9  genetic relatedness between individuals derived from genotypic similarity to further

10  decompose random effects into kinship effects between individuals. In its original application

11  to imputing missing phenotypes, PHENIX outperforms other imputation approaches when

12  the heritability ($h^2$) of a phenotype increased [9]. Similarly, when randomly removing and re-

13  imputing 10% of observed data, the quality of the imputation of PD clinical assessments was

14  in general better when considering the genetic relatedness between individuals as compared

15  to excluding this information (**Supplementary Figure 7**), suggesting that the resulting

16  phenotypic axes better capture PD heterogeneity when including genetic information.

17  Moreover, we found a higher agreement between the phenotypic axes derived by integrating

18  the genetic relationship between patients of different cohorts than when the phenotypic axes

19  were derived ignoring the genetic relationships (**Supplementary Figure 8**). Specifically, the

20  coefficient of determination reflecting the agreement between the axes derived from the

21  Discovery and those derived from the PPMI cohorts were from Axis 1 to 3: 0.665 (p=0.048),

22  0.914 (p=0.003) and 0.754 (p=0.025) when including the genetic similarity between patients

23  as compared to 0.604 (p=0.069), 0.908 (p=0.003) and 0.001 (p=0.991) without. Together,

24  these findings demonstrate that the integration of genetic relationship between patients

25  enhances the resulting phenotypic axes' ability to reproducibly capture PD clinical variation.

1    **Metanalysis of Genome Wide Association Studies with phenotypic axes as unique and**

2    **universal quantitative traits**

3          Each phenotypic axis provides a quantitative trait enabling the genetics underlying

4    patient variation to be studied by performing a Genome Wide Association Study (GWAS) via

5    a regression model with the covariates age, gender, and two genetic principal components (to

6    account for any underlying population substructure) in each individual cohort. As three

7    phenotypic axes were similar across each individual cohort (*Discovery*, *Tracking* and PPMI)

8    and to increase statistical power to detect an significant association, we conducted a meta-

9    analysis of each phenotypic axis genome-wide association studies using a common set of

10   4211937 variants across 3088 individuals. A significant departure from the expected

11   quantiles was observed for Axis 1 (meta-analysis combining the summary statistic of three

12   individual GWAS [*Discovery*-Tracking-PPMI]) (**Supplementary Figure 9**), but no variant

13   surpassed genome-wide significance (**Supplementary Figure 10**). Although we did not

14   observe a significant genome-wide association, the use of universal phenotypic axes

15   significantly unable us to conduct meta-analysis and thus to increase the statistical power to

16   identify genetic variants through their ability to align differently deeply phenotyped cohorts

17   and reduce the number of traits tested.

18         Next, we re-examined genetic associations for each of the three phenotypic axes for

19   three major PD risk genes, namely *SNCA*, *GBA* and *LRRK2*. We found a indicative local

20   association signal but however un-significant at GWA level with Phenotypic Axis 1 for a

21   variant in *SNCA*: 4: 90758437  (p-value=$1.7 \times 10^{-4}$, **Supplementary Figure 11A**) which is in

22   high LD with rs1348224 ($r^2 > 0.8$), a SNP previously associated with PD with dementia and

23   dementia with Lewy bodies [20]. SNP rs1348224:G allele (minor allele) had a negative effect

24   on Phenotypic Axis 1, thus a protective effect for cognitive impairment, which is consistent

25   with a protective effect for PD with dementia  and dementia with Lewy bodies previously

1    reported for this locus [20]. We also found a indicative local association signal (p-

2    value=$1.1 \times 10^{-4}$) with Phenotypic Axis 3 for an intronic variant in *LRRK2* (**Supplementary**

3    **Figure 11B).** Both *SNCA* and *LRRK2* variants were each nominally associated with only one

4    phenotypic axis (**Supplementary Table 4**), suggesting distinct pathogenic mechanisms.

5    **Parkinson patients carrying a high genetic risk for Alzheimer's are more**

6    **likely to develop a more aggressive form of Parkinson's**

7    To better understand the genetics risk factors influencing the phenotypic axis, we

8    calculated a disease-risk relatedness matrix in the MPMM, based on genetic variants

9    associated with different complex human traits. For example, replacing the overall genetic

10   similarity by how similar people are in their risk of diabetes or depression. By examining the

11   proportion of phenotypic variation explained by different phenotypic axis derived using these

12   different disease risks as compared to the original phenotypic axes derived using the entire

13   genotype, we showed that the phenotypic axis derived using Alzheimer's disease (AD)

14   genetic risk significantly outperforms (capture more patient phenotypic variation) the original

15   phenotypic axes (**Fig. 4)**

16   This result proposes that PD patients carrying a high genetic risk for AD are more likely to

17   develop a more aggressive form of PD including dementia symptoms: Axis 1 represents

18   worsening non-tremor motor phenotypes, anxiety and depression accompanied by a decline

19   of the cognitive function (preprint Table 1 & Fig. 2) Testing this hypothesis in the PPMI

20   cohort, we found a significant relationship between Phenotypic Axis 1 and CSF Aβ1-42, a

21   biomarker strongly associated with future conversion to dementia. Secondly, as for AD

22   genetics risk and unlikely to PD genetics risk, we found an association of phenotypic axis 1

23   risk variants with microglia-expressed genes, in both the SN and the cortex suggesting that

24   the neuro-inflammation play a key role in the  development a more aggressive form of PD,

25   but not in the PD onset-risk. Finally we observed the  phenotypic axis one is associated with

1    rapid progression of multiple clinical symptoms suggesting that AD genetic risk score in PD

2    patients could be used as a predictor of progression

## 1    Discussion

2    We propose here a novel approach to quantifying diverse patient phenotypes on a

3    continuous scale via the use of phenotype axes. This approach overcomes many of the

4    limitations associated with the clustering methods previously used to classify PD

5    heterogeneity. By applying our approach to three independent and deeply phenotyped

6    cohorts, we demonstrate the universality of these axes of phenotypic variation amongst PD

7    patients. We also showed that our axes are robustly derived when reducing the number of

8    clinical features considered and, unlike other dimensionality reduction methods, the PHENIX

9    MPMM approach is the only method tested here that is able to identify the same phenotypic

10   axes underlying PD patient variation between individuals from different cohorts. The

11   phenotypic axes have multiple applications in PD precision medicine. We found that PD

12   patients carry on a high genetic risk load for Alzheimer's disease can develop a more clinical

13   aggressive PD form including dementia symptoms.

14   Our approach was able to identify representative quantitative variables that are

15   clinically relevant to previously-defined categorical PD subtypes. A number of known

16   comorbidities were represented among the phenotype axes. Anxiety and depression are

17   highly correlated in PD patients, both of which are correlated with Axes 1 and 2 [26]. Rigidity

18   and bradykinesia are also linked, possibly due to shared physiology [27], and varied in the same

19   direction along Axis 3. Lawton *et al.* reported five PD subgroups, by using the same

20   *Discovery* cohort but following a k-means clustering approach [6]. We examined the

21   distribution of phenotypic axis score across these five PD subgroups (**Supplementary Figure**

22   **16**) and noted that the 5[th] subgroup of patients, characterised by severe motor, non-motor and

23   cognitive disease, with poor psychological well-being clinical symptoms, were systematically

24   associated with high severity score for all three of our phenotypic axes. Inversely, the first PD

25   subgroup characterised by mild motor and non-motor disease (group affected by fewer

1    clinical symptoms) displayed a low severity score for our three phenotypic axes.

2    Furthermore, we observed that the individuals of subgroups 4 and 5, characterised by poor

3    psychological well-being, had high severity scores for phenotypic axis 2, the axis most

4    associated with depression and anxiety symptoms. These observations demonstrate some

5    consistency between subgroups defined with k-means and our phenotypic axis severity score.

6    The agreement of these phenotype axes with previously observed correlations provides

7    further support for underlying biological themes, but their reinterpretation as robust

8    continuous traits likely provides a better approximation of how the underlying biology

9    contributes, as opposed to a cut-off off for a phenotype. Specifically, the unimodal character

10   of the phenotypic axis distributions (**Supplementary Figure 17**) suggests here that the

11   development of continuous measures is more appropriate than clustering according to an

12   arbitrary threshold.

13        The phenotypic axes identified were robust in terms of the number of clinical features

14   considered and enable the alignment of patients from different cohorts with different clinical

15   phenotyping structures. The corollary is that Phenix did not require the variable selection

16   common in PD clustering approaches, and it can also guide clinicians in determining which

17   clinical assessments are essential to capture PD heterogeneity. Deep phenotyping is

18   burdensome to both patient and clinician and many of the measures exploited here are

19   compound scores summarising aspects of functioning. Further work identifying the

20   minimally burdensome observations that enable robust scoring of patients along these

21   phenotypic axes would facilitate their utility and adoption across the PD clinical community,

22   bringing increased power to the discovery of influencing factors. Finally, the MPMM

23   approach can be readily extended to include longitudinal data to determine the phenotypic

24   axes associated with disease progression while simultaneously dealing with missing data,

25   which is a common problem in longitudinal studies.

1    In conclusion, these universal axes have the potential to accelerate our understanding

2    of how PD presents in individual patients, providing more robust and objective quantitative

3    traits through which patients may be appropriately compared, through which the underlying

4    disease-modifying mechanism can be understood and appropriately stratified/personalised

5    therapeutic strategies and treatments can be developed.

6

## Acknowledgments

22

## Conflict of interest

24   The authors declare that they have no competing interests.

# References

1.  Foltynie T, Brayne C, Barker RA. The heterogeneity of idiopathic Parkinson's disease. *J Neurol* 2002; **249**(2)**:** 138-145.

2.  Malek N, Swallow DM, Grosset KA, Lawton MA, Marrinan SL, Lehn AC *et al.* Tracking Parkinson's: Study Design and Baseline Patient Data. *J Parkinsons Dis* 2015; **5**(4)**:** 947-959.

3.  PPMI. PPMI.

4.  Szewczyk-Krolikowski K, Tomlinson P, Nithi K, Wade-Martins R, Talbot K, Ben-Shlomo Y *et al.* The influence of age and gender on motor and non-motor features of early Parkinson's disease: initial findings from the Oxford Parkinson Disease Center (OPDC) discovery cohort. *Parkinsonism Relat Disord* 2014; **20**(1)**:** 99-105.

5.  Fereshtehnejad SM, Zeighami Y, Dagher A, Postuma RB. Clinical criteria for subtyping Parkinson's disease: biomarkers and longitudinal progression. *Brain* 2017; **140**(7)**:** 1959-1976.

6.  Lawton M, Baig F, Rolinski M, Ruffman C, Nithi K, May MT *et al.* Parkinson's Disease Subtypes in the Oxford Parkinson Disease Centre (OPDC) Discovery Cohort. *J Parkinsons Dis* 2015; **5**(2)**:** 269-279.

7.  Lawton M, Ben-Shlomo Y, May MT, Baig F, Barber TR, Klein JC *et al.* Developing and validating Parkinson's disease subtypes and their motor and cognitive progression. *Journal of Neurology, Neurosurgery & Psychiatry* 2018.

8.  Erro R, Picillo M, Vitale C, Palladino R, Amboni M, Moccia M *et al.* Clinical clusters and dopaminergic dysfunction in de-novo Parkinson disease. *Parkinsonism Relat Disord* 2016; **28:** 137-140.

9.  Dahl A, Iotchkova V, Baud A, Johansson A, Gyllensten U, Soranzo N *et al.* A multiple-phenotype imputation method for genetic studies. *Nat Genet* 2016; **48**(4)**:** 466-472.

10. Malek N, Weil RS, Bresner C, Lawton MA, Grosset KA, Tan M *et al.* Features of GBA-associated Parkinson's disease at presentation in the UK Tracking Parkinson's study. *J Neurol Neurosurg Psychiatry* 2018.

11. Parkinson Progression Marker I. The Parkinson Progression Marker Initiative (PPMI). *Prog Neurobiol* 2011; **95**(4)**:** 629-635.

12. Parkes M, Cortes A, van Heel DA, Brown MA. Genetic insights into common pathways and complex relationships among immune-mediated diseases. *Nat Rev Genet* 2013; **14**(9)**:** 661-673.

13. Nalls MA, Bras J, Hernandez DG, Keller MF, Majounie E, Renton AE *et al.* NeuroX, a fast and efficient genotyping platform for investigation of neurodegenerative diseases. *Neurobiol Aging* 2015; **36**(3): 1605 e1607-1612.

14. Nalls MA, Keller MF, Hernandez DG, Chen L, Stone DJ, Singleton AB *et al.* Baseline genetic associations in the Parkinson's Progression Markers Initiative (PPMI). *Mov Disord* 2016; **31**(1): 79-85.

15. Purcell S, Neale B, Todd-Brown K, Thomas L, Ferreira MA, Bender D *et al.* PLINK: a tool set for whole-genome association and population-based linkage analyses. *Am J Hum Genet* 2007; **81**(3): 559-575.

16. Andreassen OA, Zuber V, Thompson WK, Schork AJ, Bettella F, Consortium P *et al.* Shared common variants in prostate cancer and blood lipids. *Int J Epidemiol* 2014; **43**(4): 1205-1214.

17. Andreassen OA, Djurovic S, Thompson WK, Schork AJ, Kendler KS, O'Donovan MC *et al.* Improved detection of common variants associated with schizophrenia by leveraging pleiotropy with cardiovascular-disease risk factors. *Am J Hum Genet* 2013; **92**(2): 197-209.

18. Zuber V, Jonsson EG, Frei O, Witoelar A, Thompson WK, Schork AJ *et al.* Identification of shared genetic variants between schizophrenia and lung cancer. *Sci Rep* 2018; **8**(1): 674.

19. Winsvold BS, Bettella F, Witoelar A, Anttila V, Gormley P, Kurth T *et al.* Shared genetic risk between migraine and coronary artery disease: A genome-wide analysis of common variants. *PLoS One* 2017; **12**(9): e0185663.

20. Guella I, Evans DM, Szu-Tu C, Nosova E, Bortnick SF, Group SCS *et al.* alpha-synuclein genetic variability: A biomarker for dementia in Parkinson disease. *Ann Neurol* 2016; **79**(6): 991-999.

21. Andreassen OA, Thompson WK, Schork AJ, Ripke S, Mattingsdal M, Kelsoe JR *et al.* Improved detection of common variants associated with schizophrenia and bipolar disorder using pleiotropy-informed conditional false discovery rate. *PLoS Genet* 2013; **9**(4): e1003455.

22. Noyce AJ, Kia DA, Hemani G, Nicolas A, Price TR, De Pablo-Fernandez E *et al.* Estimating the causal influence of body mass index on risk of Parkinson disease: A Mendelian randomisation study. *PLoS Med* 2017; **14**(6): e1002314.

23. Brainstorm C, Anttila V, Bulik-Sullivan B, Finucane HK, Walters RK, Bras J *et al.* Analysis of shared heritability in common disorders of the brain. *Science* 2018; **360**(6395).

24. Nalls MA, Saad M, Noyce AJ, Keller MF, Schrag A, Bestwick JP *et al.* Genetic comorbidities in Parkinson's disease. *Hum Mol Genet* 2014; **23**(3): 831-841.

1   25.   Birtwistle J, Baldwin D. Role of dopamine in schizophrenia and Parkinson's disease.
2         *Br J Nurs* 1998; **7**(14)**:** 832-834, 836, 838-841.
3
4   26.   Menza MA, Robertson-Hoffman DE, Bonapace AS. Parkinson's disease and anxiety:
5         comorbidity with depression. *Biol Psychiatry* 1993; **34**(7)**:** 465-470.
6
7   27.   Berardelli A, Rothwell JC, Thompson PD, Hallett M. Pathophysiology of
8         bradykinesia in Parkinson's disease. *Brain* 2001; **124**(Pt 11)**:** 2131-2146.
9
10
11

1    **Figure Legends**

2    **Fig. 1**. **The clinical phenotypes of two independent deeply phenotyped Parkinson's**

3        **disease cohorts identify the same phenotypic axes.** Results were consistent in two

4        independents cohorts (842 *Discovery* and 1807 *Tracking* patients). Examination of

5        these two separate Parkinson's disease cohorts, using independent derivation of the

6        phenotypic axes in each, showed significant correlations between each cohort's first

7        three axes. Correlations between the axes from each cohort are Axis 1 r=0.92 (p=3 x

8        10-13), Axis 2 r=0.89 (p=4 x 10-11), Axis 3 r=0.72 (p=5 x 10-6). The correlation

9        coefficient (x-axis) between each axis derived in each cohort (blue: *Discovery* vs red:

10       *Tracking*) and each clinical observation (y-axis) is shown.

11

12   **Fig. 2. The reduced dimensions in other dimensionality reduction methods fail to align**

13       **between differently but deeply phenotyped UK and US Parkinson's disease**

14       **cohorts.** We compared the ability of different dimensionality reduction methods

15       (independent component analysis (ICA), Multidimensional scaling (MDS), Principal

16       component analysis (PCA) and phenotypic axis based on the PHENIX multiple

17       phenotype mixed model) to phenotypically align two deeply phenotyped Parkinson's

18       disease cohorts, specifically the *Discovery* (842 individuals) and PPMI (439 sporadic

19       Parkinson's disease) cohorts. The x-axis and y-axis represent the correlation

20       coefficient between each continuous variable with clinical observation associated with

21       a specific symptom category in *Discovery* and PPMI cohort respectively. Each

22       column panel and colour of points ("Axis") represents the dimension level of each

23       underlying dimension.  All points on the diagonal would represent a perfect

24       phenotypic alignment of both cohorts. We examined the relationship between

1      correlation derived from both cohorts by performing a linear regression: R^2 and p

2      correspond to the coefficient of determination and the p-value respectively. Only the

3      dimensions discovered by the MPMM model, PHENIX, show a significant

4      relationship between both cohorts: MPMM phenotypic axes ($R^2$=0.86, p=2x10-8),

5      MDS ($R^2$=0.11, p=0.18), ICA ($R^2$=0.17, p=0.16) and PCA ($R^2$=0.31, p=0.06).

6

7

8      **Fig. 3.     The correlation of individual clinically-measured Parkinson's disease**

9      **phenotypes with an underlying Phenotypic Axis 1.** Modelling patient clinical

10      phenotypes as a combination of genetic and environmental factors revealed three

11      phenotypic severity axes (**Fig.1**), each representing a continuous pattern of variation

12      between multiple co-varying clinical phenotypes. In Axis 1 (shown), (A) clinical

13      measures relating to anxiety and depression and apathy are significantly and

14      positively correlated with an individual's score along this axis; patients with a higher

15      axis score have more severe mood and neuropsychiatric problems. (B) The severity of

16      motor phenotypes is positively correlated with this phenotypic axis; patients with a

17      higher axis score is associated with more severe motor phenotypes (C) Cognitive tests

18      were negatively correlated with this component (the patients that score high in these

19      cognitive tests have less cognitive impairments); individuals with a high score for this

20      component suffer from more severe anxiety, depression and displayed more cognitive

21      impairment and motor symptoms.

22

1

2    **Fig. 4. Alzheimer's disease phenotypic axis significantly outperforms the original**

3         **phenotypic axis.** These heatmap plots represents for two first phenotypic axes

4         (left=V1 and right=V2), in the 3 cohorts (row1=Discovery, row2=Tracking & row 3=

5         ) the excess(red) or deficit of the phenotypic variance explained by the different

6         disease axes (columns) compared with the original phenotypic axes.

7

8

1     **Table1:** Correlation between each axis and each clinical phenotypic measure

| Category | Clinical Observation | | | Axis1 | Axis2 | Axis3 |
|---|---|---|---|---|---|---|
| Behavior | BDI total | Measure of the depression | + | 0.60 | 0.60 | 0.01 |
| | Leeds Anxiety Total | Measure of the anxiety | + | 0.51 | 0.55 | 0.05 |
| | Leeds Depression | Measure of the depression | + | 0.51 | 0.62 | 0.00 |
| | QUIP all | Impulsive-Compulsive Disorders | + | 0.12 | 0.24 | 0.03 |
| | UPDRS apathy | Apathy | + | 0.40 | 0.39 | 0.18 |
| | UPDRS fatigue | Fatigue | + | 0.49 | 0.49 | 0.08 |
| | UPDRS hallucinations | Hallucinations | + | 0.17 | 0.17 | -0.02 |
| Autonomic | Constipation | Quantitative measure of constipation | + | -0.15 | -0.07 | 0.01 |
| | Orthostatic | Blood pression from sitting/lying to stand up | + | 0.17 | -0.09 | -0.06 |
| | UPDRS constipation | Constipation | + | 0.38 | 0.33 | -0.08 |
| | UPDRS pain | Pain | + | 0.47 | 0.47 | -0.01 |
| Cognitive | Education years | Number of years of education | - | -0.21 | -0.23 | 0.16 |
| | MMSE total | Measure of cognitive ability | - | -0.27 | -0.07 | 0.17 |
| | MOCA total | Measure of cognitive ability | - | -0.31 | -0.06 | 0.23 |
| | Phonemic fluency | Number of words beginning with a particular letter | - | -0.26 | 0.03 | 0.14 |
| | Sementic fluency | Number of animals and the number of boy names | - | -0.28 | 0.09 | 0.15 |
| | BMI | Body Mass index | + | 0.16 | 0.09 | -0.08 |
| | CGIC | Clinical global impression of change | + | 0.05 | -0.07 | 0.08 |
| | Disease Duration | Disease Duration | + | 0.24 | 0.19 | -0.07 |
| Motors | Flamingo time | Time that a person can stand on one leg | - | -0.46 | -0.03 | 0.16 |
| | Getgo average | Time taken for an individual to get up from a chair, walk three meters, turn around, walk back to the chair and sit down. | + | 0.52 | 0.04 | -0.16 |
| | Purdue assembly | Test to measure manual dexterity | - | -0.37 | 0.16 | 0.10 |
| | Purdue total | Test to measure manual dexterity | - | -0.41 | 0.18 | 0.09 |
| | UPDRS arms | Arms | + | 0.63 | -0.50 | 0.78 |
| | UPDRS bradykinesia | Bradykinesia | + | 0.63 | -0.40 | 0.57 |
| | UPDRS faceneck | Face/neck problems | + | 0.26 | -0.22 | 0.12 |
| | UPDRS I | Non Motor Aspects of Experiences of Daily Living | + | 0.68 | 0.67 | 0.02 |
| | UPDRS II | Motor Aspects of Experiences of Daily Living | + | 0.76 | 0.30 | 0.05 |
| | UPDRS III | Motors Examination | + | 0.71 | -0.46 | 0.61 |
| | UPDRS IV | Motors complications | + | 0.16 | 0.16 | 0.05 |
| | UPDRS laterality | Unilateral | + | -0.03 | -0.01 | 0.15 |
| | UPDRS legs | Legs | + | 0.59 | -0.31 | 0.44 |
| | UPDRS postural | Postural | + | 0.64 | -0.02 | -0.09 |
| | UPDRS rigidity | Rigidity | + | 0.51 | -0.27 | 0.35 |

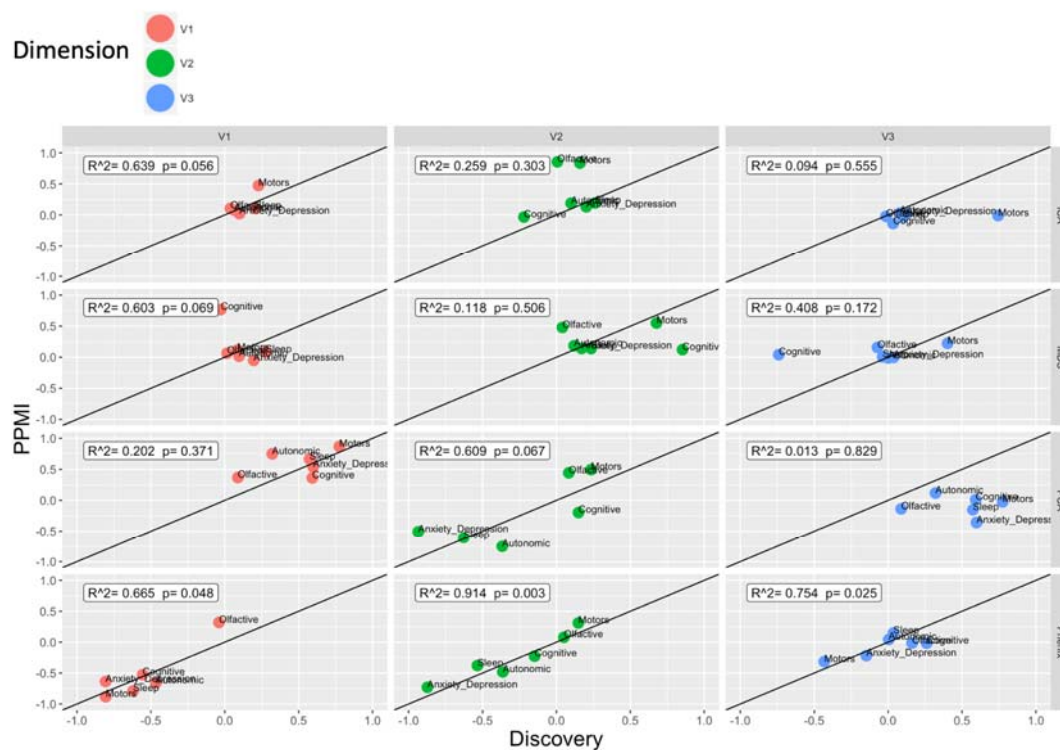| | | | | | | |
|---|---|---|---|---|---|---|
| | UPDRS speech | Speech | + | 0.22 | -0.07 | -0.04 |
| | UPDRS tremor | Tremor | + | 0.20 | -0.40 | 0.58 |
| Sleep | ESS total | Measure of daytime sleepiness | + | 0.31 | 0.22 | -0.07 |
| | RBD total | Measure of REM Sleep behavior disorder | + | 0.29 | 0.29 | -0.03 |
| e | | | | | | |
| | Sniff total | Smell identifications | - | 0.01 | 0.06 | 0.12 |
| Drug | LEDD total | Quantitative measure of the amount of Parkinson's disease medication | + | 0.31 | 0.27 | -0.22 |

(1)  A high score for a clinical measure indicates **more (+)** or **less (-)** issue for the patient.

(2)  The correlation coefficient under and above |0.25| are indicated in gray or blue/red respectively

(3)  Red and blue cells indicates when a high phenotypic axis score are associated with more and less clinical issues for the patient respectively.
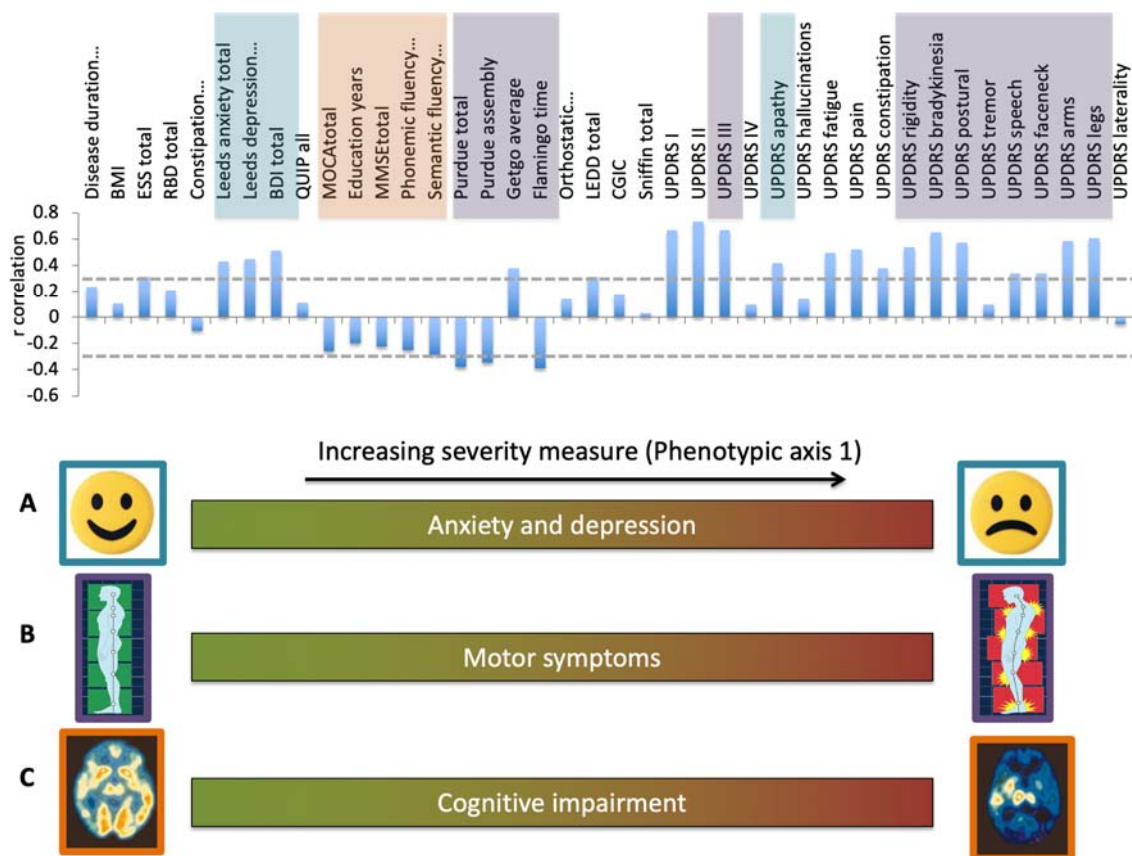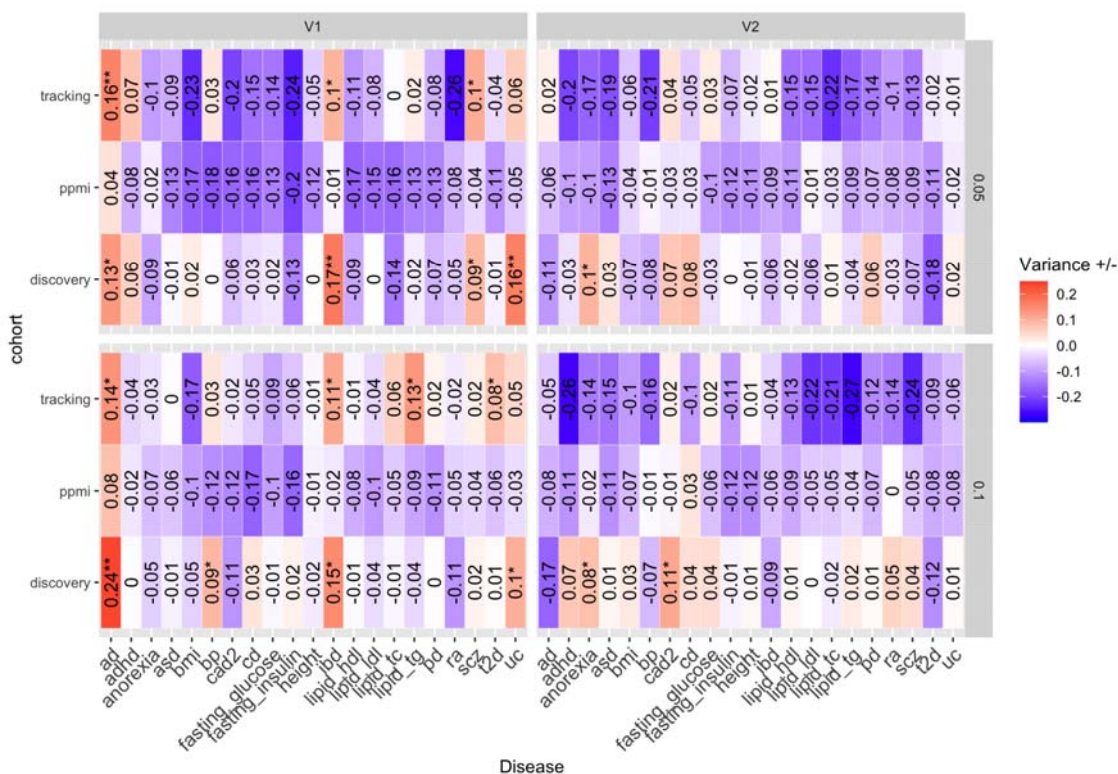
1    **Fig.1**
2



3
4

1    **Fig. 2**
2



3

1   **Fig.3**
2



3
4
5
6
7
8

1    **Fig. 4**

2



3
4
5