

Title

Integration of proteomics with genomics and transcriptomics increases the diagnostic rate of Mendelian disorders

Authors

Robert Kopajtich^{†1,2}, **Dmitrii Smirnov**^{†1,2}, **Sarah L. Stenton**^{†1,2}, Stefan Loipfinger³, Chen Meng⁴, Ines F. Scheller⁵, Peter Freisinger⁶, Robert Baski⁷, Riccardo Berutti^{1,2}, Jürgen Behr⁴, Martina Bucher⁸, Felix Distelmaier⁹, Mirjana Gusic^{1,2}, Maja Hempel¹⁰, Lea Kulterer^{1,2}, Johannes Mayr¹¹, Thomas Meitinger¹, Christian Mertes³, Metodi D. Metodiev¹², Agnieszka Nadel^{1,2}, Alessia Nasca^{13,14}, Akira Ohtake¹⁵, Yasushi Okazaki¹⁶, Rikke Olsen¹⁷, Dorota Piekutowska-Abramczuk¹⁸, Agnès Rötig¹², René Santer¹⁹, Detlev Schindler²⁰, Abdelhamid Slama²¹, Christian Staufner²², Tim Strom¹, Patrick Verloo²³, Jürgen-Christoph von Kleist-Retzow²⁴, Saskia B. Wortmann^{11,25}, Vicente A. Yépez^{1,3}, Costanza Lamperti¹³, Daniele Ghezzi^{13,14}, Kei Murayama²⁶, Christina Ludwig⁴, Julien Gagneur^{3,27}, Holger Prokisch^{*1,2}

† These authors contributed equally to this work

* Corresponding author

Affiliations

1. Institute of Human Genetics, Technical University of Munich, Munich, Germany
2. Institute of Neurogenomics, Helmholtz Zentrum München, Munich, Germany
3. Department of Informatics, Technical University of Munich, Garching, Germany

NOTE: This preprint reports new research that has not been certified by peer review and should not be used to guide clinical practice.

4. Bavarian Center for Biomolecular Mass Spectrometry (BayBioMS), Technical University of Munich, Freising, Germany
5. Institute of Computational Biology, Helmholtz Zentrum München, Munich, Germany
6. Department of Pediatrics, Klinikum Reutlingen, Reutlingen, Germany
7. Leeds Teaching Hospitals, NHS Trust, Leeds, UK
8. Institut für Klinische Genetik, Olgahospital, Stuttgart, Germany
9. Department of General Pediatrics, Neonatology and Pediatric Cardiology, Heinrich-Heine-University, Düsseldorf, Germany
10. Institute of Human Genetics, University Medical Center Hamburg-Eppendorf, Hamburg, Germany
11. Department of Pediatrics, Salzburger Landeskliniken and Paracelsus Medical University, Salzburg, Austria
12. UMR1163, Université Paris Descartes, Sorbonne Paris Cité, Institut Imagine, Paris, France
13. Unit of Medical Genetics and Neurogenetics, Fondazione IRCCS Istituto Neurologico Carlo Besta, Milan, Italy
14. Department of Pathophysiology and Transplantation, University of Milan, Milan, Italy
15. Center for Intractable Diseases, Department of Pediatrics & Clinical Genomics, Faculty of Medicine, Saitama Medical University, Saitama, Japan
16. Diagnostics and Therapeutic of Intractable Diseases, Intractable Disease Research Center, Graduate School of Medicine, Juntendo University, Tokyo, Japan
17. Research Unit for Molecular Medicine, Department for Clinical Medicine, Aarhus University and Aarhus University Hospital, Aarhus, Denmark
18. Department of Medical Genetics, Children's Memorial Health Institute (CMHI) Warsaw, Warsaw, Poland

19. Department of Pediatrics, University Medical Center Eppendorf, Hamburg, Germany
20. Department of Human Genetics, University of Würzburg, Würzburg, Germany
21. Department of Biochemistry, Reference Center for Mitochondrial Disease, FILNEMUS, Bicêtre University Hospital, University of Paris-Saclay, Assistance Publique-Hôpitaux de Paris, Le Kremlin-Bicêtre, France
22. Division of Neuropediatrics and Pediatric Metabolic Medicine, Center for Child and Adolescent Medicine, University Hospital Heidelberg, Heidelberg, Germany
23. Division of Child Neurology and Metabolism, Department of Pediatrics, Ghent University Hospital, Ghent, Belgium
24. University of Cologne, Faculty of Medicine and University Hospital Cologne, Department of Pediatrics, Cologne, Germany
25. Radboud Centre for Mitochondrial Diseases (RCMM), Amalia Children's Hospital, Radboudumc, Nijmegen, The Netherlands
26. Center for Medical Genetics, Department of Metabolism, Chiba Children's Hospital, Chiba, Japan
27. Institute of Computational Biology, Helmholtz Zentrum München, Neuherberg, Germany

Corresponding author contact information

Dr. Holger Prokisch

Institut für Humangenetik, Klinikum rechts der Isar, Technische Universität München

Trogerstraße 32, 81675 München, Germany

E-mail: prokisch@helmholtz-muenchen.de

Number: +49 89 3187 2890

Manuscript

By lack of functional evidence, genome-based diagnostic rates cap at approximately 50% across diverse Mendelian diseases. Here we demonstrate the effectiveness of combining genomics, transcriptomics, and, for the first time, proteomics and phenotypic descriptors, in a systematic diagnostic approach to discover the genetic cause of mitochondrial diseases. On fibroblast cell lines from 145 individuals, tandem mass tag labelled proteomics detected approximately 8,000 proteins per sample and covered over 50% of all Mendelian disease-associated genes. By providing independent functional evidence, aberrant protein expression analysis allowed validation of candidate protein-destabilising variants and of variants leading to aberrant RNA expression. Overall, our integrative computational workflow led to genetic resolution for 21% of 121 genetically unsolved cases and to the discovery of two novel disease genes. With increasing democratization of high-throughput omics assays, our approach and code provide a blueprint for implementing multi-omics based Mendelian disease diagnostics in routine clinical practice.

The current ACMG recommendation for interpretation of genetic variants (Richards et al., 2015) attaches high importance to functional validation in designation of a variant as pathogenic or likely pathogenic. For this reason, systematic application of RNA sequencing (RNA-seq) has proven valuable in reducing the diagnostic shortfall of whole exome sequencing (WES) or whole genome sequencing (WGS) by providing a molecular diagnosis to 10% of unsolved cases with a mitochondrial disease (Kremer et al., 2017) and up to 35% in other disease cohorts (Cummings et al., 2017, Gonorazky et al., 2019, Fresard et al., 2019). However, while proteomics has been used to validate variants of uncertain significance (VUS) in single

cases (Kremer et al., 2017, Lake et al., 2017, Borno et al., 2019, Stojanovski et al., 2020), the utility of systematic application into a diagnostic pipeline has yet to be explored.

Within the GENOMIT project (genomit.eu) we had analysed approximately 1,000 clinically suspected mitochondrial disease cases by WES/WGS and gathered corresponding Human Phenotype Ontology (HPO) terms for automated phenotype integration. Mitochondrial diseases are a prime example of the diagnostic challenge faced in human genetics given their vast clinical and genetic heterogeneity. In-keeping with previous studies (Stenton and Prokisch 2020), we reached a diagnosis by WES/WGS analysis for approximately 50% of the cases. Here, selecting 143 of these mitochondrial disease cases (121 unsolved and 22 solved positive controls) plus two healthy controls with available fibroblast cell lines, we performed RNA-seq and quantitative tandem mass tag (TMT) labelled proteomics in an integrative multi-omic approach (Fig.1a, Supplementary Fig. 1a-c) (see online Methods).

With the detection of approximately 12,000 transcripts and 8,000 proteins per sample (Supplementary Fig. 1d), a median of 91% (n=353) and 80% (n=310) of mitochondrial disease gene products, and 59% (n=2535) and 51% (n=2159) of all Mendelian disease gene products were quantified per sample in RNA-seq and proteomics, respectively, deeming fibroblasts an easily accessible tissue with high disease gene coverage and an excellent resource for the study of mitochondrial diseases (Supplementary Fig. 1e).

To identify genes with aberrant RNA expression, we performed three outlier analyses, i) aberrant expression levels, ii) aberrant splicing, and iii) monoallelic expression of rare variants via the DROP pipeline (Yépez et al., 2021). To identify aberrant protein expression, we developed the algorithm PROTRIDER which estimates deviations from expected protein

intensities while controlling for known and unknown sources of proteome-wide variation (see online Methods). In 18 positive controls with nuclear-encoded variants, detection of protein expression outliers in 14 (77.8%) validated our proteomic approach (Fig. 1c). In four positive controls with mtDNA-encoded variants, there was no significant change in protein expression resulting in a total validation rate of 64% across 22 positive controls (Supplementary Table 1).

In our cohort of 121 unsolved cases we first investigated those with variants prioritised in the WES/WGS analysis, spanning a total of 26 unique alleles across 21 cases, and mostly missense in nature (Supplementary Table 2). Variant pathogenicity was validated in 14 cases (67%) by nominally significant protein underexpression (Fig. 1d), of which five were also validated by aberrant RNA expression (Supplementary Fig. 1f). Moreover, proteomics was valuable in rejecting the prioritised variant, such as in the mitochondrial targeting sequence of *MRPL53* which associated with normal expression of both *MRPL53* and the large mitoribosomal subunit (Supplementary Fig. 1g).

Our matched genome, transcriptome, and proteome datasets together with protein expression outlier calls allowed us to investigate how aberrant RNA and protein expression relate to one another in the context of rare genetic variation, to our knowledge for the first time. After multiple-testing correction, we identified a median of two aberrantly expressed transcripts and six aberrantly expressed proteins per sample (Supplementary Fig. 2a). Though less than half of the RNA outliers resulted in significant protein outliers, possibly due to buffering mechanisms on the protein level, artefact, or lack of power, the majority (77%) did result in protein outliers (Fig 2a). The expression outliers were stratified into three classes: RNA-only, protein-only, and RNA-and-protein outliers. We focused on the two thirds of outliers that are underexpressed, as evidence for impaired function. All three classes of underexpression

outliers were significantly enriched for rare variants in their encoding gene (Supplementary Fig. 2b-c). Within RNA-only outliers there was enrichment for splice, stop, and frameshift variants, in line with RNA expression outlier studies (Li et al., 2017). In contrast, protein-only outliers captured the functional consequence of missense variants and in-frame indels, demonstrating significant enrichment for coding variants. Substantially more RNA-and-protein outliers (15%) could be explained by potentially biallelic rare variants, compared to RNA-only and protein-only outliers (approximately 5%, respectively). An additional 25% of RNA-and-protein outliers were associated with rare heterozygous variants. Protein outliers without rare variants in the encoding gene may be explained indirectly as a consequence of protein complex instability due to a defect in one of the interaction partners (Kremer et al., 2017, Lake et al., 2017, Borno et al., 2019). Collectively, these genome-wide observations emphasize the complementarity of proteomics to RNA-seq in capturing the functional impact of rare genetic variation. Moreover it shows the sensitivity of the approach not only to biallelic variation, a hallmark of a recessive inheritance mode, but also to mono-allelic variation, i.e. those responsible for dominant diseases. If these biallelic or single variants are in-keeping with the inheritance mode and phenotype of the known disease gene, the detection of an outlier may lead to diagnosis of the patient.

Aiming to pinpoint pathogenic genes and variants for those remaining cases without prioritized VUS from WES/WGS, we next combined aberrant expression analysis with patient phenotypic annotations (Fig. 2a, Supplementary Fig. 2d-e). Focussing on significant underexpression outliers (median 4 per sample), a median of one outlier matched with the patient phenotype as described by HPO annotations (see online Methods). Manual inspection and clinical interpretation of the outliers resulted in the diagnosis of 12 cases (11%) (Fig. 2b-e, Supplementary Fig. 3). In four cases, the identified protein-only outlier led to the diagnosis by

providing functional validation of VUS in a gene not previously prioritised, yet in-keeping with the phenotype and mode of inheritance of the disease (Fig. 2b). In eight cases, we identified the diagnosis as a significant RNA-and-protein outlier. Of these eight cases, one case had a one exon deletion in *NFUI* identified by follow-up WGS in compound heterozygosity with a missense variant resulting in 21% residual protein (z-score -7.8). One case has a heterozygous missense variant in *MORC2* resulting in 69% residual protein (z-score -4.5) (Fig. 2c). Four cases demonstrated aberrant splicing resulting in protein outliers, such as a homozygous near splice variant in *MRPL44* (z-score -5.5), deep intronic variants in *TIMMDC1* (z-score -6.2) and *MRPS25* (z-score -4.8), and in one case a direct splice variant on one allele and a unique combination of two frequent intronic variants on the second allele (allele frequency 7.2% and 21.8%, respectively) causing exon skipping in *DARS2* (z-score -6.3) (Fig. 2d). In one case, compound heterozygous variants in *VPS11* originally prioritised by WES were not validated given normal protein expression. However, underexpression of *MRPS25* and five subunits of the small mitoribosomal subunit led to the diagnosis, exemplifying the added value of proteomics in detecting the consequence on all detected proteins and complexes (Fig. 2e). Finally, in two cases, our integrated omics approach led to the identification of novel mitochondrial disease genes, *MRPL38* and *LIG3*. The *MRPL38* outlier (z-score -5.8) (Fig. 3a) illuminated a pathogenic 5'UTR deletion (Fig. 3b). The functional relevance was confirmed by reduced abundance of the large mitoribosomal subunit (Fig. 3c) resulting in a severe reduction in mitochondrial translation rate rescued by the re-expression of wild-type *MRPL38* (Fig. 3d). The *LIG3* outlier (z-score -4.2) (Fig. 3e) reprioritised a heterozygous nonsense variant within the mitochondrial targeting sequence (Fig. 3f) affecting only the mitochondrial isoform *in trans* with a deep intronic variant causing aberrant splicing (Fig. 3f). As a dual localized nuclear and mitochondrial DNA ligase, a defect in *LIG3* was expected to impact mitochondrial DNA replication, supported by mtDNA depletion and a combined OXPHOS defect in the muscle

biopsy (Fig. 3g), significantly decreased protein levels of mtDNA encoded gene products (Fig. 3h), and impaired mtDNA repopulation (Fig. 3i). The downstream functional consequence of the *LIG3* variants was reflected by 63 additional protein outliers. The four cases solved by protein-only outliers, included a hemizygous X-linked *NDUFB11* missense variant resulting in aberrant protein underexpression (z-score -4.1) and pathologically low abundance of respiratory chain complex (RCC) I (50%) with no rare variants within any other RCCI subunit. The reduction in RCCI was most pronounced in the ND4-module to which *NDUFB11* belongs (44% remaining, lowest in dataset), in-keeping with a second confirmed *NDUFB11* case (55% remaining, second lowest in dataset). The detection of this variant in the unaffected grandfather indicated incomplete penetrance. Attributing pathogenicity to variants of incomplete penetrance, even in the presence of a phenotypic match, is an outstanding challenge in human genetics. However, in cases where reduced activity is causative of disease, proteomics has the power to classify variants affecting protein complex abundance. This was also demonstrated for a homozygous variant in *DNAJC30* in two cases, by providing evidence for the loss-of-function character of an incompletely penetrant missense variant, as recently reported in a cohort of 27 families (Stenton et al., 2021).

To summarise, leveraging on advanced proteomics we quantified a substantial fraction of expressed proteins, determined their normal physiological range in fibroblasts, and called protein outliers in a robust manner. By developing an integrated multi-omic analysis pipeline, we establish a clinical decision support tool for the diagnosis of Mendelian disorders. The power of proteomics is demonstrated by validation and detection of the molecular diagnosis in 26 of 121 (21%) unsolved WES/WGS cases, of which in 11 (42%) we detect downstream functional evidence on the complex level, explaining in total more than 100 outliers in these 26 cases (Fig. 4). Our code is freely available

(<https://prokischlab.github.io/omicsDiagnostics/>). An interactive web interface allows the user to browse all results and could serve as a basis for developing future integrative multi-omics diagnostic interfaces. We used TMT-labelling, a proteomics technique quantifying the very same peptides for all samples of a batch. This greatly facilitates detection of under-expression outliers compared to conventional untargeted mass-spectrometry which suffers from widespread missing values in low intensity ranges. Though RNA-seq did not allow interpretation of missense variants, it provided independent cumulative evidence and guided the identification of causative splice variants in half of the solved cases. Moreover, RNA-seq has a deeper coverage of expressed genes, capturing 50% more genes. It therefore remains useful for lowly expressed proteins. To identify novel diagnoses we applied stringent significance filtering (FDR<0.1) and focussed on underexpression outliers with a phenotype match, leading to one protein outlier per sample in median. However, with the integration of multiple levels of omics information and phenotype descriptors, relaxed significance thresholds may in future be considered. Our approach depends on an available tissue, encouraging clinicians to be proactive and opportunistic in biosampling, specifically when follow-up visits are unlikely. Given the increasing democratization of proteomics we envisage its implementation in clinical practice to advance diagnostics by routine integration of functional data.

online Methods

Study cohort

All individuals included in the study or their legal guardians provided written informed consent before evaluation, in agreement with the Declaration of Helsinki and approved by the ethical committees of the centres participating in this study, where biological samples were obtained.

All studies were completed according to local approval of the ethical committee of the Technical University of Munich.

Cell culture

Primary fibroblast cell lines were cultured as per Kremer et al., 2017.

Whole exome sequencing (WES)

Whole exome sequencing was performed as per Kremer et al., 2017. SAMtools v.0.1.19 and GATK v.4.0 and called on the targeted exons and regions from the enrichment kit with a +/- 50bp extension.

Variant annotation and handling

Variant Effect Predictor (McLaren et al., 2016) from Ensembl (Zerbino et al., 2018) was used to annotate genetic variants with minor allele frequencies from the 1000 Genomes Project (1000 Genome Consortium, 2015), gnomAD (Karczewski et al., 2020), and the UK Biobank (Bycroft et al., 2018), location, deleteriousness scores and predicted consequence with the highest impact among all possible transcripts. Variants with minor allele frequency less than 1% across all cohorts were considered as rare. Genes harbouring one rare allele were classified as rare, with two or more rare alleles - potentially biallelic. ACMG guidelines for variant classification were implemented with the InterVar software (Li and Wang, 2015).

Gene-phenotypic matching

Phenotype similarity was calculated as symmetric semantic similarity score with R::PCAN package (Godard and Page, 2016). We considered genes to match phenotypically if the symmetric semantic similarity between the gene and the case HPO annotations was larger or

equal to 2 (Köhler et al., 2009; Frésard et al., 2019) (Supplementary Fig. 2d). Affected organ systems were visualized with R:: gganatogram (Maag 2018), based on patient's HPO phenotypes corresponding to the third level of HPO ontology (Köhler et al., 2019).

RNA-sequencing

Non-strand specific RNA-seq was performed as per Kremer et al., 2017. Strand-specific RNA-seq was performed according to the TruSeq Stranded mRNA Sample Prep LS Protocol (Illumina, San Diego, CA, USA). Processing of RNA sequencing files was performed as per Kremer et al., 2017.

Detection of aberrant RNA expression, aberrant splicing, and mono-allelic expression

RNA-seq analysis was performed using DROP (Yepez et al., 2021), an integrative workflow that integrates quality controls, expression outlier calling with OUTRIDER (Brechtmann et al., 2018), splicing outlier calling with FRASER (Mertes et al., 2020), and mono-allelic expression with a negative binomial test (Kremer et al., 2017). We used as reference genome the GRCh37 primary assembly, release 29, of the GENCODE project (Frankish et al., 2019) which contains 60,829 genes. RNA expression outliers were defined as those with a false-discovery rate ≤ 0.1 . Splicing outliers were defined as those with a gene-level false-discovery rate ≤ 0.1 and a deviation of the observed percent-spliced-in or splicing efficiency from their expected value larger than 0.3. Mono-allelic expression was assessed only for heterozygous single nucleotide variants reported by WES analysis. We retained mono-allelic expression calls at a false discovery rate ≤ 0.05 . Aberrant events of all three types were further inspected using the Integrative Genome Viewer (Robinson et al., 2011).

Mass spectrometric sample preparation

Proteomics was performed at the BayBioMS core facility at the Technical University Munich, Freising, Germany. Fibroblast cell pellets containing 0.5 million cells were lysed under denaturing conditions in urea containing buffer and quantified using BCA Protein Assay Kit (Thermo Scientific). 15 µg of protein extract were further reduced, alkylated and the tryptic digest was performed using Trypsin Gold (Promega). Digests were acidified, desalted and TMT-labeling was performed according to (Zecha et al., 2019) using TMT 10-plex labelling reagent (Thermo Fisher Scientific). Each TMT-batch consisted of 8 patient samples and 2 reference samples common to all batches to allow for data normalization between batches. Each TMT 10-plex peptide mix was fractionated using trimodal mixed-mode chromatography as described (Yu et al., 2017). LC-MS measurements were conducted on a Fusion Lumos Tribrid mass spectrometer (Thermo Fisher Scientific) which was operated in data-dependent acquisition mode and multi-notch MS3 mode. Peptide identification was performed using MaxQuant version 1.6.3.4 (Tyanova et al., 2016) and protein groups obtained. Missing values were imputed with the minimal value across the dataset.

Transcriptome-proteome matching

In order to determine the correct assignment of proteome and transcriptome assay from the same sample, we correlated the gene counts with the protein intensities (Supplementary Fig. 4). The spearman ranked correlation test was applied to all transcriptome-proteome combinations, using the `cor.test` function from R. The distribution of the correlation values are plotted and in the case of mismatch two distinctive populations will appear. Correlations greater than 0.2 correspond to matching samples. Only protein intensities greater than 10,000 and genes with at least 50 counts were considered. Protein intensities were log-transformed and centered. RNA counts were normalized by sequencing depth using size factors (Love et al.,

2014), log-transformed and centered. The 2,000 genes with the highest dispersion (as computed by OTRIDER) were selected.

Detection of aberrant protein expression with PROTRIDER

To detect protein expression outliers while controlling for known and unknown sources of proteome-wide variations, we employed a denoising autoencoder based method, analogous to methods for calling RNA expression outliers (Brechtmann et al., 2018) and splicing outliers (Mertes et al., 2020). Specifically, sizefactor normalized and log-transformed protein intensities were centred protein-wise and used as input to a denoising autoencoder model with three layers (encoder, hidden space, decoder). As protein intensities varied strongly between batches, we included the batch as a covariate in the input of the encoder and in the input of the decoder (Supplementary Fig. 5a-b). For a given encoding dimension q , we fit the autoencoder by minimizing the mean squared error loss over the non-missing data. The optimal encoding dimension of the autoencoder was determined by artificially injecting outliers and selecting the encoding dimension that yielded the best area under the precision-recall curve (AUPRC) of recovering these injected outliers. For this dimension fitting procedure, artificial outliers were generated with a frequency of 1 per 1000. An outlier log-transformed intensity $x_{i,j}^0$ for a sample i and a protein j was generated by shifting the observed log-transformed intensity $x_{i,j}$ by $z_{i,j}$ times the standard deviation σ_j of $x_{i,j}$, with the absolute value of $z_{i,j}$ being drawn from a log-normal distribution with the mean of the logarithm equal to 3 and the standard deviation of the logarithm equal to 1.6, and with the sign of $z_{i,j}$ either up or down, drawn uniformly:

$$x_{i,j}^0 = x_{i,j} + z_{i,j} \cdot \sigma_j.$$

After the autoencoder model was fit to the data, statistical testing of the observed log-transformed intensities $x_{i,j}$ with respect to the expected log-transformed intensities $\mu_{i,j}$ modelled by the autoencoder was performed, using two-sided Gaussian p-values $p_{i,j}$ for sample i and protein j defined as

$$p_{i,j} = 2 \cdot \min\{N(x_{i,j}|\mu_{i,j}, \sigma^{\text{res}}_j), 1 - N(x_{i,j}|\mu_{i,j}, \sigma^{\text{res}}_j)\},$$

where σ^{res}_j is the protein-wise standard deviation of the autoencoder residuals $x_{i,j} - \mu_{i,j}$. Finally, p-values were corrected for multiple testing per sample with the method of Benjamini and Yekutieli (Benjamini and Yekutieli, 2001). During the entire process of fitting the autoencoder model as well as the statistical tests, missing data was masked as unavailable and ignored. We refer to this method as PROTRIDER in the following.

Benchmark of PROTRIDER against limma

As no method for outlier detection in proteomics data was established yet, we benchmarked our method against an approach that is based on limma (Smyth 2005), which was developed for differential expression analyses on microarray data and assesses statistical significance with a moderated t-statistic. We used recalibrated protein data which has been adjusted with respect to the two identical control samples in each MS-run as the input for limma and included the sex, batch and instrument annotation to adjust for confounding factors. To be able to use limma for outlier detection, we tested each sample against all other samples. We evaluated the performance of both methods based on precision-recall curves of detecting the known category I defects (Supplementary Fig. 5c-f). In this benchmark, PROTRIDER showed superior performance, as it was able to recover more known defects while reporting fewer outliers in

total (median per sample of 6 vs. 8 for the limma based approach). Therefore, we decided to adopt PROTRIDER for the detection of aberrant protein expression.

Enrichment of genetic variants in outlier genes

We focused our analysis only on the genes where both RNA-and-protein levels were quantified, per every sample and limited it to the genes that were detected as outliers at least once in our cohort. Variants were stratified into six classes according to their impact on the protein sequence, defined by a combination of VEP (McLaren et al., 2016) annotations as follows: Stop (stop_lost, stop_gained), splice (splice_region_variant, splice_acceptor_variant, splice_donor_variant), frameshift (frameshift_variant), coding (missense_variant, protein_altering_variant, inframe_insertion, inframe_deletion), synonymous (synonymous_variant, stop_retained_variant) and non-coding (3_prime_UTR_variant, 5_prime_UTR_variant, downstream_gene_variant, upstream_gene_variant, intron_variant, non_coding_transcript_exon_variant, mature_miRNA_variant, intron_variant, intergenic_variant, regulatory_region_variant). Enrichment analysis was performed similarly as described by Li et al 2017, by modelling with logistic regression of each outlier category (RNA only, protein only, RNA-and-protein over- or underexpression) as a function of standardized variant class. For each gene, detected as an outlier of a particular category, the remaining set of individuals served as controls. Proportions of outlier genes were calculated by assignment of one variant class (out of six) with the highest significant enrichment in the corresponding outlier category.

Detection of aberrantly expressed protein complexes

Detection of aberrantly expressed protein complexes was performed similar to the differential protein complex expression method described by Zhou et al., 2019. Specifically, the quantified

proteins were mapped to the protein complex database CORUM (v3.0) (Giurgiu et al., 2019) or to the mitochondria-related subset of HGNC gene groups by gene names. We considered the protein complexes of four subunits or more and with at least 50% of the subunits quantified. For each sample i and protein complex k , we computed $y_{i,k}$, the mean deviation of observed versus expected protein intensities across all detected subunits (expressed in \log_2 fold-change and as estimated by PROTRIDER or LIMMA). For each protein complex k , we fitted by maximum likelihood a Gaussian on all $y_{i,k}$ with mean μ_k and standard deviation σ_k using the `fitdistr` function from the R package MASS (Venables and Ripley 2002). The two-sided Gaussian p-values for sample i and protein complex k was then computed as:

$$p_{i,k} = 2 * \min \{ N(y_{i,k} | \mu_k, \sigma_k), 1 - N(y_{i,k} | \mu_k, \sigma_k) \},$$

To correct the p-values for multiple testing, the method of Benjamini and Yekutieli (Benjamini and Yekutieli, 2001) was applied per every sample.

Mitochondrial translation assays

Metabolic labelling of mitochondrial proteins was performed essentially as described previously (Ruzzenente et al., 2018). In brief, fibroblasts derived from individuals #102875 and 98521 were incubated in methionine- and cysteine-free DMEM medium supplemented with 10% dialyzed FBS, GlutaMAX, sodium pyruvate (ThermoFisher Scientific, Montigny-le-Bretonneux, France), 100 mg/ml emetine dihydrochloride to block cytosolic protein synthesis and 400 μCi EasyTag EXPRESS35S Protein Labelling Mix (PerkinElmer, Villebon-sur-Yvette, France). Labelling was performed for 30 min followed by a further incubation for 10 min in standard growth medium. Equal amounts of total cell lysates were fractionated by SDS-PAGE and newly synthesized proteins were quantified by autoradiography.

Data and code availability

The proteomic raw data and MaxQuant search files have been deposited to the ProteomeXchange Consortium (<http://proteomecentral.proteomexchange.org>) via the PRIDE partner repository and can be accessed using the dataset identifier PXD022803. Code to reproduce the analysis is available via GitHub at github.com/prokischlab/omicsDiagnostics/.

Online resources

Code to reproduce the figures: <https://github.com/prokischlab/omicsDiagnostics/tree/master>

Web interfaces: <https://prokischlab.github.io/omicsDiagnostics/#readme.html>

PRIDE: <https://www.ebi.ac.uk/pride/archive/projects/PXD022803>

DROP: <https://github.com/gagneurlab/drop>

GTEEx Portal: <https://www.gtexportal.org/home/>

OMIM database: www.omim.org

CORUM: <https://mips.helmholtz-muenchen.de/corum/>

HGNC: <https://www.genenames.org>

Acknowledgements

This study was supported by a German Federal Ministry of Education and Research (BMBF, Bonn, Germany) grant to the German Network for Mitochondrial Disorders (mitoNET, 01GM1906D), the German BMBF and Horizon2020 through the E-Rare project GENOMIT (01GM1920A), the ERA PerMed project PerMiM (01KU2016A), the German BMBF through the e:Med Networking funds AbCD-Net (FKZ 01ZX1706A), the German Research Foundation/Deutsche Forschungsgemeinschaft (DI 1731/2-2), the AFM-Telethon grant (#19876), the Practical Research Project for Rare/Intractable Diseases from the Japan Agency

for Medical Research and Development, AMED (JP19ek0109273, JP20ek0109468, JP20kk0305015, JP20ek0109485), a CMHI grant (S145/16), a PMU-FFF grant (A-20/01/040-WOS), the Pierfranco and Luisa Mariani Foundation (CM23), and the Italian Ministry of Health (GR-2016-02361494, GR-2016-02361241). We would like to thank Caterina Terrile, Franziska Hackbarth, and Hermine Kienberger for their excellent laboratory assistance as well as Miriam Abele for mass spectrometric support at the BayBioMS. We thank the “Cell line and DNA Bank of Genetic Movement Disorders and Mitochondrial Diseases” of the Telethon Network of Genetic Biobanks (grant GTB12001J) and Eurobiobank Network which supplied biological specimens.

Author Contributions

Conceived and supervised the study, H.P; performed experiments, R.K, L.K, C.Lu, D.G, and M.M; analyzed and interpreted results, D.S, S.L, I.S, C.M, V.Y, D.G, M.M, R.K, S.L.S, J.G, H.P; provided essential materials, all authors; wrote the manuscript, S.L.S, H.P, J.G, R.K, and D.S; edited manuscript, all authors.

Competing Interests Statement

A.O. declares a consigned research fund (SBI Pharmaceuticals Co., Ltd.). All other authors declare no conflict of interest.

References

- 1000 Genomes Consortium, G. P., Auton, A., Brooks, L. D., Durbin, R. M., Garrison, E. P., & Kang, H. M. (2015). A global reference for human genetic variation. *Nature*, 526(7571), 68-74.
- Benjamini, Y., & Yekutieli, D. (2001). The control of the false discovery rate in multiple testing under dependency. *Annals of statistics*, 1165-1188.
- Borna, N. N., Kishita, Y., Kohda, M., Lim, S. C., Shimura, M., Wu, Y., ... & Okazaki, Y. (2019). Mitochondrial ribosomal protein PTCD3 mutations cause oxidative phosphorylation defects with Leigh syndrome. *neurogenetics*, 20(1), 9-25.
- Brechtmann, F., Mertes, C., Matusevičiūtė, A., Yépez, V. A., Avsec, Ž., Herzog, M., ... & Gagneur, J. (2018). OUTRIDER: A statistical method for detecting aberrantly expressed genes in RNA sequencing data. *The American Journal of Human Genetics*, 103(6), 907-917.
- Bycroft, C., Freeman, C., Petkova, D., Band, G., Elliott, L. T., Sharp, K., ... & Marchini, J. (2018). The UK Biobank resource with deep phenotyping and genomic data. *Nature*, 562(7726), 203-209.
- Cummings, B. B., Marshall, J. L., Tukiainen, T., Lek, M., Donkervoort, S., Foley, A. R., ... & Estrella, E. (2017). Improving genetic diagnosis in Mendelian disease with transcriptome sequencing. *Science translational medicine*, 9(386).
- Frankish, A., Diekhans, M., Ferreira, A. M., Johnson, R., Jungreis, I., Loveland, J., ... & Flicek, P. (2019). GENCODE reference annotation for the human and mouse genomes. *Nucleic acids research*, 47(D1), D766-D773.
- Frésard, L., Smail, C., Ferraro, N. M., Teran, N. A., Li, X., Smith, K. S., ... & Balliu, B. (2019). Identification of rare-disease genes using blood transcriptome sequencing and large control cohorts. *Nature medicine*, 25(6), 911-919.

- Giurgiu, M., Reinhard, J., Brauner, B., Dunger-Kaltenbach, I., Fobo, G., Frishman, G., ... & Ruepp, A. (2019). CORUM: the comprehensive resource of mammalian protein complexes—2019. *Nucleic acids research*, 47(D1), D559-D563.
- Godard, P., & Page, M. (2016). PCAN: phenotype consensus analysis to support disease-gene association. *BMC bioinformatics*, 17(1), 1-9.
- Gonorazky, H. D., Naumenko, S., Ramani, A. K., Nelakuditi, V., Mashouri, P., Wang, P., ... & Mathews, K. D. (2019). Expanding the boundaries of RNA sequencing as a diagnostic tool for rare mendelian disease. *The American Journal of Human Genetics*, 104(3), 466-483.
- Karczewski, K. J., Francioli, L. C., Tiao, G., Cummings, B. B., Alföldi, J., Wang, Q., ... & MacArthur, D. G. (2020). The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature*, 581(7809), 434-443.
- Köhler, S., Carmody, L., Vasilevsky, N., Jacobsen, J. O. B., Danis, D., Gouridine, J. P., ... & Robinson, P. N. (2019). Expansion of the Human Phenotype Ontology (HPO) knowledge base and resources. *Nucleic acids research*, 47(D1), D1018-D1027.
- Köhler, S., Schulz, M. H., Krawitz, P., Bauer, S., Dölken, S., Ott, C. E., ... & Robinson, P. N. (2009). Clinical diagnostics in human genetics with semantic similarity searches in ontologies. *The American Journal of Human Genetics*, 85(4), 457-464.
- Kremer, L. S., Bader, D. M., Mertes, C., Kopajtich, R., Pichler, G., Iuso, A., ... & Koňářiková, E. (2017). Genetic diagnosis of Mendelian disorders via RNA sequencing. *Nature communications*, 8(1), 1-11.
- Lake, N. J., Webb, B. D., Stroud, D. A., Richman, T. R., Ruzzenente, B., Compton, A. G., ... & Thorburn, D. R. (2017). Biallelic mutations in MRPS34 lead to instability of the small mitoribosomal subunit and Leigh syndrome. *The American Journal of Human Genetics*, 101(2), 239-254.

- Li, Q., & Wang, K. (2017). InterVar: clinical interpretation of genetic variants by the 2015 ACMG-AMP guidelines. *The American Journal of Human Genetics*, 100(2), 267-280.
- Li, X., Kim, Y., Tsang, E. K., Davis, J. R., Damani, F. N., Chiang, C., ... & Li, A. (2017). The impact of rare variation on gene expression across tissues. *Nature*, 550(7675), 239-243.
- Love, M. I., Huber, W., & Anders, S. (2014). Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome biology*, 15(12), 1-21.
- Maag, J. L. (2018). gganatogram: An R package for modular visualisation of anatograms and tissues based on ggplot2. *F1000Research*, 7.
- McLaren, W., Gil, L., Hunt, S. E., Riat, H. S., Ritchie, G. R., Thormann, A., ... & Cunningham, F. (2016). The ensembl variant effect predictor. *Genome biology*, 17(1), 1-14.
- Mertes, C., Scheller, I. F., Yépez, V. A., Çelik, M. H., Liang, Y., Kremer, L. S., ... & Gagneur, J. (2021). Detection of aberrant splicing events in RNA-seq data using FRASER. *Nature communications*, 12(1), 1-13.
- Reyes Tellez, A., Melchionda, L., Nasca, A., Carrara, F., Lamantea, E., Zanolini, A., ... & Bonato, S. (2015). RNASEH1 Mutations Impair mtDNA Replication and Cause Adult-Onset Mitochondrial Encephalomyopathy.
- Richards, S., Aziz, N., Bale, S., Bick, D., Das, S., Gastier-Foster, J., ... & Voelkerding, K. (2015). Standards and guidelines for the interpretation of sequence variants: a joint consensus recommendation of the American College of Medical Genetics and Genomics and the Association for Molecular Pathology. *Genetics in medicine*, 17(5), 405-423.
- Robinson, J. T., Thorvaldsdóttir, H., Winckler, W., Guttman, M., Lander, E. S., Getz, G., & Mesirov, J. P. (2011). Integrative genomics viewer. *Nature biotechnology*, 29(1), 24-26.
- Ruzzenente, B., Assouline, Z., Barcia, G., Rio, M., Boddaert, N., Munnich, A., ... & Metodiev, M. D. (2018). Inhibition of mitochondrial translation in fibroblasts from a patient expressing

the KARS p.(Pro228Leu) variant and presenting with sensorineural deafness, developmental delay, and lactic acidosis. *Human mutation*, 39(12), 2047-2059.

Schlieben, L. D., & Prokisch, H. (2020). The dimensions of primary mitochondrial disorders. *Frontiers in Cell and Developmental Biology*, 8.

Smyth, G. K. (2005). Limma: linear models for microarray data. In *Bioinformatics and computational biology solutions using R and Bioconductor* (pp. 397-420). Springer, New York, NY.

Stenton, S. L., & Prokisch, H. (2020). Genetics of mitochondrial diseases: Identifying mutations to help diagnosis. *EBioMedicine*, 56, 102784.

Stenton, S. L., Sheremet, N. L., Catarino, C. B., Andreeva, N., Assouline, Z., Barboni, P., ... & Prokisch, H. (2021). Impaired complex I repair causes recessive Leber's hereditary optic neuropathy. *The Journal of Clinical Investigation*.

Stojanovski, D., Jackson, T. D., Hock, D., Palmer, C., Kang, Y., Fujihara, K., ... & Stroud, D. A. (2020). The TIM22 complex regulates mitochondrial one-carbon metabolism by mediating the import of Sideroflexins. *bioRxiv*.

Tyanova, S., Temu, T., & Cox, J. (2016). The MaxQuant computational platform for mass spectrometry-based shotgun proteomics. *Nature protocols*, 11(12), 2301.

Venables, W. N., & Ripley, B. D. (2002). *Modern applied statistics with S* fourth edition. World.

Yépez, V. A., Mertes, C., Müller, M. F., Klaproth-Andrade, D., Wachutka, L., Frésard, L., ... & Gagneur, J. (2021). Detection of aberrant gene expression events in RNA sequencing data. *Nature Protocols*, 1-21.

Yu, P., Petzoldt, S., Wilhelm, M., Zolg, D. P., Zheng, R., Sun, X., ... & Kuster, B. (2017). Trimodal mixed mode chromatography that enables efficient offline two-dimensional peptide fractionation for proteome analysis. *Analytical chemistry*, 89(17), 8884-8891.

Zecha, J., Satpathy, S., Kanashova, T., Avanesian, S. C., Kane, M. H., Clauser, K. R., ... & Kuster, B. (2019). TMT labeling for the masses: a robust and cost-efficient, in-solution labeling approach. *Molecular & cellular proteomics*, 18(7), 1468-1478.

Zerbino, D. R., Achuthan, P., Akanni, W., Amode, M. R., Barrell, D., Bhai, J., ... & Flicek, P. (2018). Ensembl 2018. *Nucleic acids research*, 46(D1), D754-D761.

Zhou, B., Yan, Y., Wang, Y., You, S., Freeman, M. R., & Yang, W. (2019). Quantitative proteomic analysis of prostate tissue specimens identifies deregulated protein complexes in primary prostate cancer. *Clinical proteomics*, 16(1), 1-18.

Figure Legends

Fig. 1: Genetic diagnosis by simultaneous genomic (WES/WGS), phenomic, transcriptomic (RNA-seq), and proteomic investigation followed by integrated analysis.

a, Multi-omic approach based on the integration of genomics (WES/WGS), transcriptomics (RNA-seq), proteomics, and phenotypic descriptors (HPO). We obtained DNA for WES/WGS from blood and RNA-seq and proteomics from fibroblasts obtained by minimally invasive skin biopsy. Functional evidence from each omic is integrated in search of a genetic diagnosis. The resultant diagnosis is thereby supported by multiple lines of robust clinical and functional evidence. Simultaneously, heterozygous and potentially biallelic genetic variants were prioritized according to their effect on the corresponding transcript(s) and protein by the identification of outliers in RNA-seq and proteomic data, in addition to aberrant splicing and monoallelic expression (MAE) of a deleterious heterozygous variants in RNA-seq data. Phenotype data complemented the analysis by gene-level prioritization based upon semantic similarity scoring. Together, omics integration allowed comprehensive gene-variant prioritization by providing insight into the effect of rare variation on expression of gene

products. An overview of our multiplexed, time-efficient, RNA-seq and proteomic sample workflow is depicted in Supplementary Fig. 1a-c. **b**, Proportion of protein-coding genes detected by RNA-seq (blue) and proteomics (red), genome-wide and among mitochondrial disease genes (Schlieben and Prokisch 2020), Neuromuscular and Neurology genes (Frésard et al., 2019), and OMIM disease genes (<https://omim.org>). **c**, Protein z-score distribution for disease-causing genes in WES/WGS solved cases (positive controls). **d**, Protein z-score distribution for disease-causing genes with prioritized variants in WES/WGS unsolved cases. In panels **c** and **d**, the points appear in red for validated cases and in green for novel cases diagnosed in our downstream systematic approach.

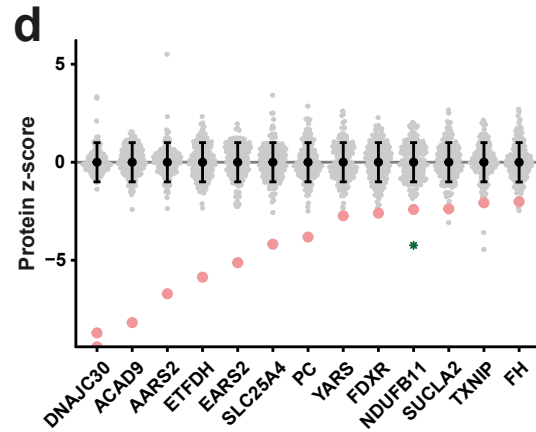
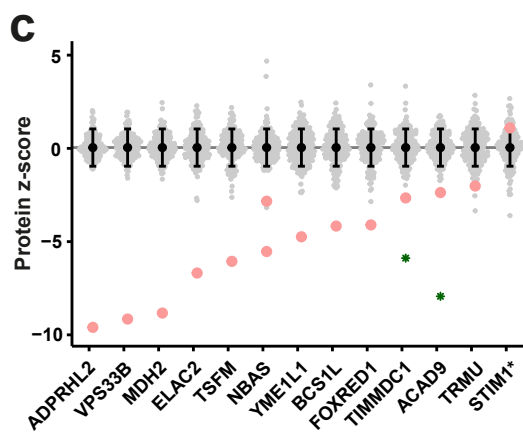
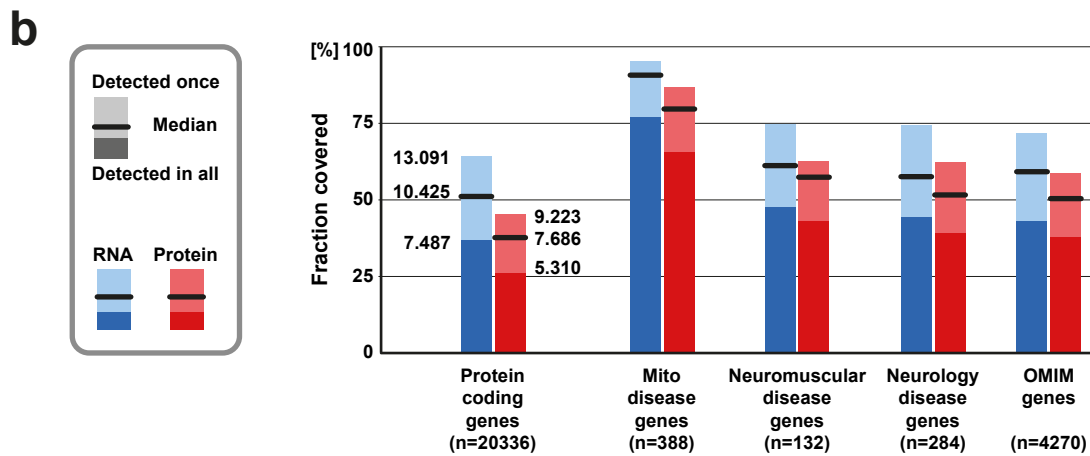
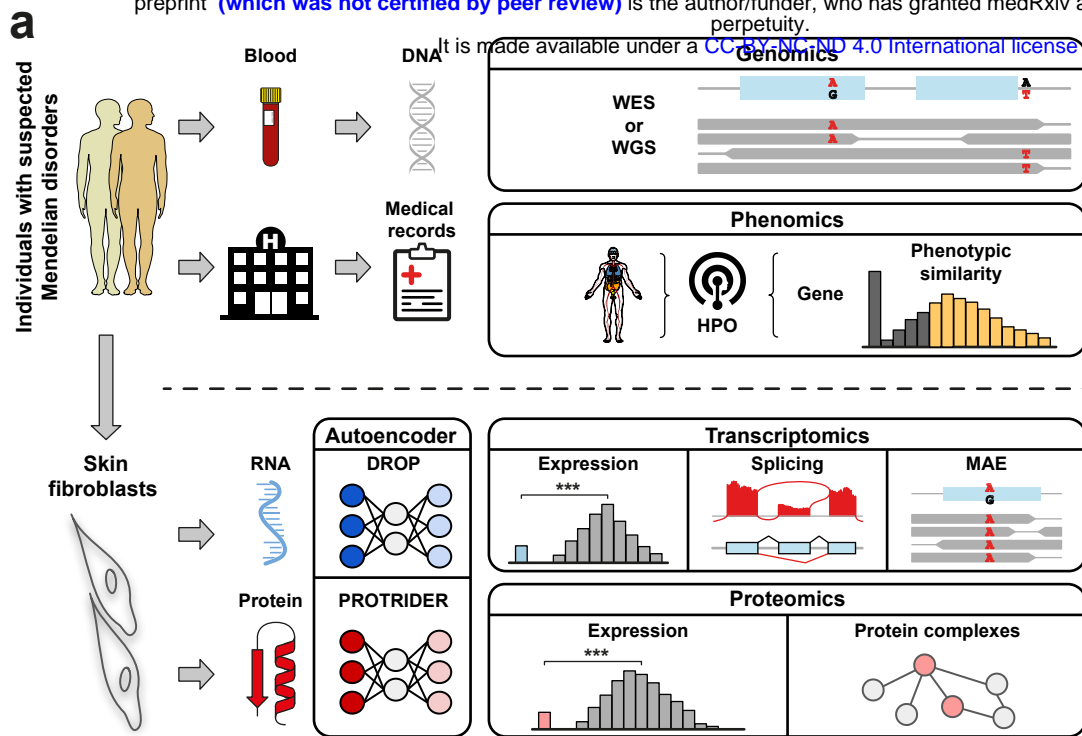
Fig. 2: The power of an integrative multi-omic approach to detect the pathogenic consequence of genetic variants. **a**, RNA z-scores (y-axis) vs. protein z-scores (x-axis) detected by RNA-seq and proteomics across all samples. The shape indicates rare variants (minor allele frequency <1%) in the encoding gene. The size indicates semantic similarity with the established disease-gene associated phenotype. The colour represents outlier class. All detected splice defects and monoallelic expression (MAE) events resulted in aberrant expression, allowing this to be used as an indicator of variant pathogenicity. **b**, Individual OM06865 presented in childhood with predominantly neurological and muscular involvement. A homozygous missense variant in the autosomal recessive disease gene *EPG5* was prioritised as a protein-only outlier. **c**, Individual OM27390 presented in infancy with failure to thrive, global developmental delay, seizures, encephalopathy, nystagmus, hypotonia, and abnormality of the basal ganglia on MRI. A heterozygous missense variant in the autosomal dominant disease gene *MORC2* was prioritized by an RNA-and-protein outlier. **g**, Individual OM75740 presented in infancy with muscular hypotonia, cardiomyopathy, and abnormalities in the cerebral white matter on MRI. *DARS2* was detected as an RNA-and-protein underexpression

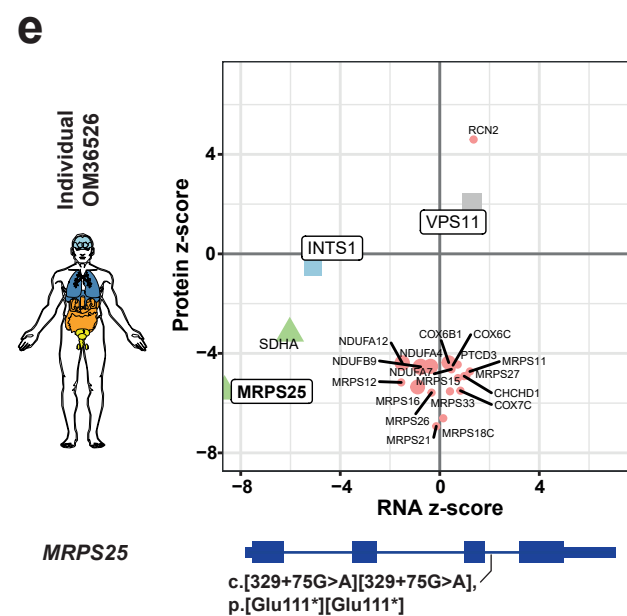
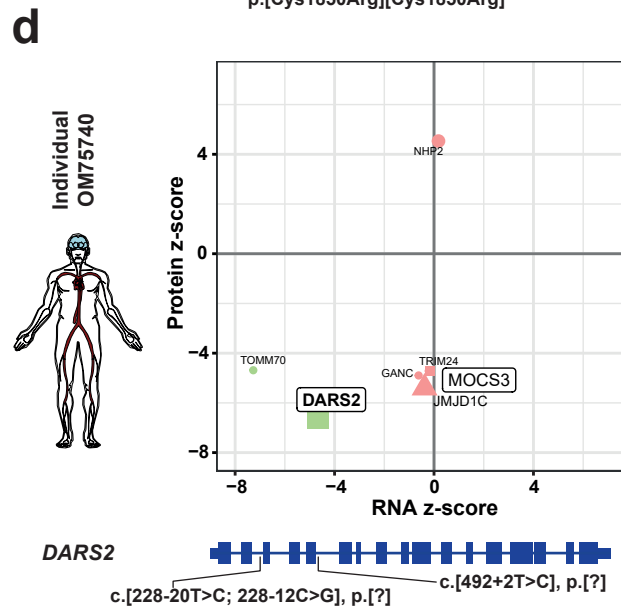
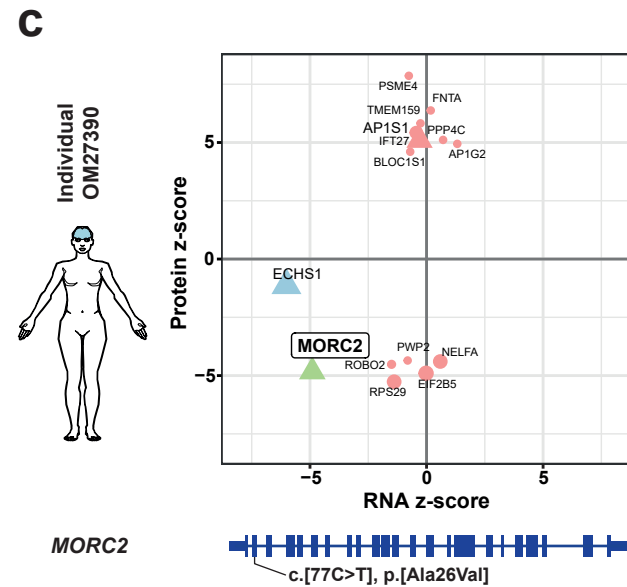
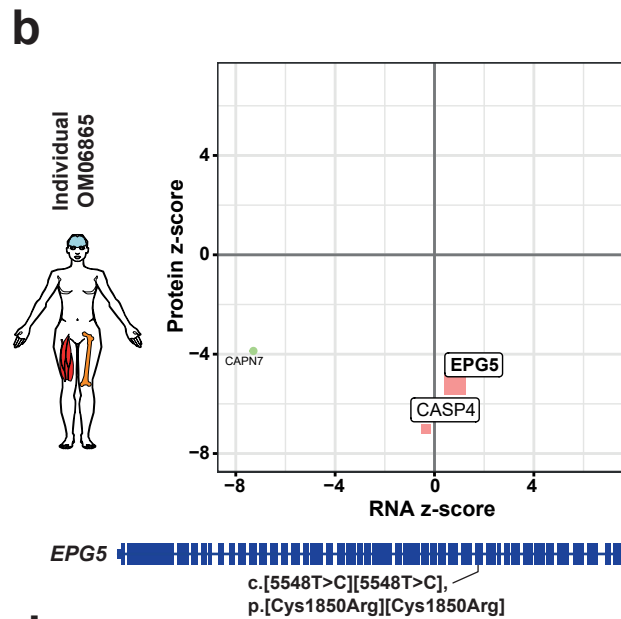
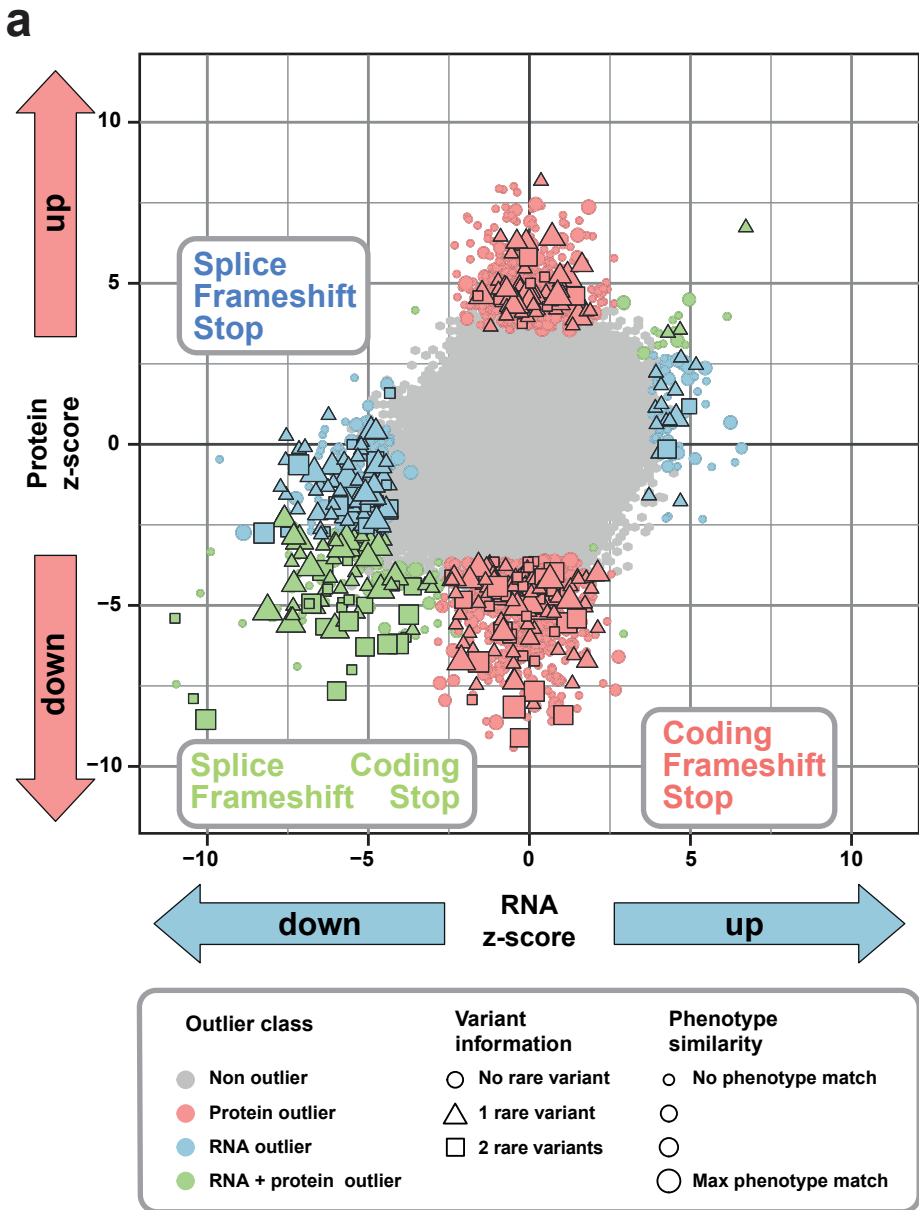
outlier explained by a combination of splice and intronic variants. **h**, Individual OM36526 presented in childhood with global developmental delay, elevated lactate, and an isolated respiratory chain complex (RCC) IV defect. The MRPS25 RNA-and-protein outlier is explained by a homozygous intronic variant demonstrating aberrant splicing. In addition, three other proteins from the small mitoribosomal subunit were underexpressed as a downstream consequence of the primary defect.

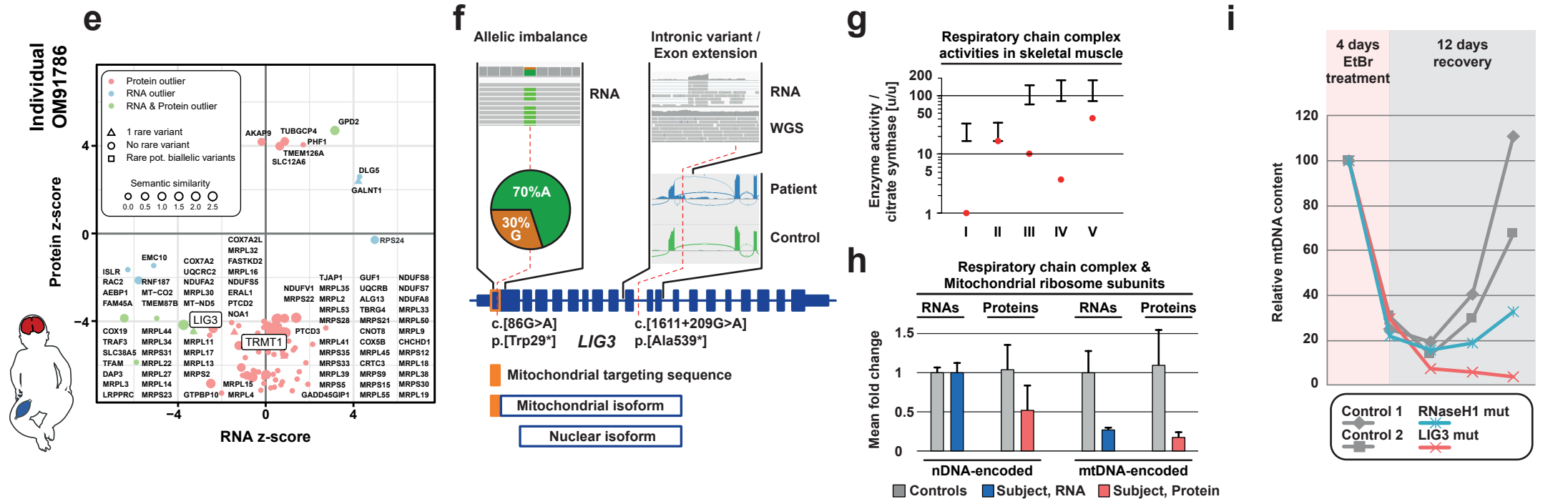
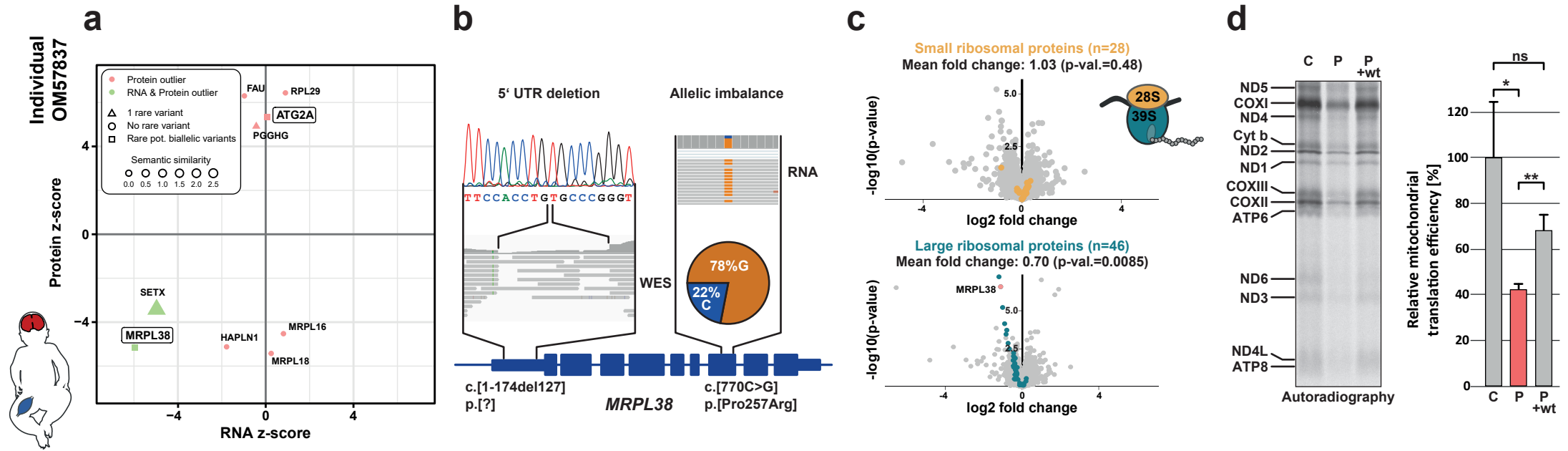
Fig. 3: Multi-omic prioritization and functional characterization of two novel mitochondrial disease genes *MRPL38* and *LIG3*. **a**, Individual OM57837 presented in infancy with global developmental delay, intellectual disability, seizures, hypotonia, symmetrical basal ganglia and brainstem abnormalities on brain MRI, and respiratory chain complex (RCC) I and IV defects. The RNA and the protein products of *MRPL38* (Mitochondrial Large Ribosomal Subunit Protein L38) were detected as underexpression outliers. **b**, A missense variant, c.[770C>G], p.[Pro257Arg], present in 78% of RNA reads indicated reduced expression of a compound heterozygous 127 bp deletion in the 5'UTR of *MRPL38*. **c**, Underexpression of *MRPL38* resulted in reduction of the large mitoribosomal subunit (n=46 detected subunits). Meanwhile, the small mitoribosomal subunit remained unchanged (n=28 detected subunits). Data are displayed as a gene-wise protein expression volcano plot of nominal ($-\log_{10}$) p-values against protein intensity \log_2 fold change. **d**, Measurement of mitochondrial translation in cultured patient-derived fibroblasts by metabolic labelling with [35S]-containing amino acids. The *MRPL38* mutant (P) showed a significantly reduced mitochondrial translation rate compared to control fibroblasts (C). **e**, Individual OM91786 presented with neonatal-onset severe encephalopathy, seizures, hypotonia, and increased serum lactate with early demise in the first weeks of life. *LIG3* was identified as an RNA-and-protein outlier. **f**, Whole-genome sequencing identified compound heterozygous

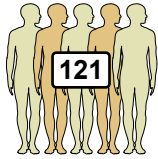
variants in *LIG3*. A nonsense variant within the mitochondrial targeting sequence affecting only the mitochondrial isoform of the protein in compound heterozygosity with a deep intronic variant leading to aberrant splicing and partial degradation of the transcript from this allele, indicated by allelic imbalance. **g**, A combined OXPHOS defect sparing nuclear encoded RCCII was present on the muscle biopsy. The black bars represent the reference range. The red point represents the patient measurement. **h**, Normal transcript level but reduced protein level of nuclear-encoded RCC subunits and ribosomal subunits (n=133) indicated complex instability. RNA-and-protein level of the 13 mtDNA encoded RNAs (11 mRNAs and 2 rRNAs) and 13 mtDNA encoded proteins were reduced. Expression is depicted as a mean-fold change compared to the mean of all other fibroblast samples. **i**, mtDNA copy number in cultured fibroblasts was investigated by qPCR during ethidium bromide induced depletion and repopulation. Impaired mtDNA repopulation was more severe than in the *RNASEH1* mutant cell line which serves as a control for a repopulation defect (Reyes et al. 2015).

Fig. 4: The value of proteomics in variant interpretation. A summary of the approach followed in the validation (left) and discovery (right) of protein outliers in our cohort of 121 unsolved mitochondrial disease cases to reach a molecular diagnosis in 26 cases (21%).





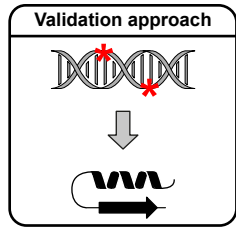




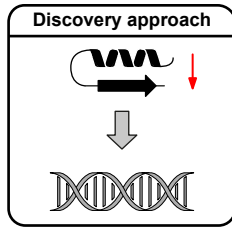
Unsolved WES/WGS cases with suspected Mendelian disease

21
with prioritized VUS

100
without prioritized VUS



7 not validated



14 diagnoses validated
(67%)

12 diagnoses discovered
(11%)

21% total diagnostic rate
(26 / 121)

42% with downstream functional evidence on the complex level
(11 / 26)

