




RESEARCH PAPER

 OPEN ACCESS 

Characterization of protein isoform diversity in human umbilical vein endothelial cells via long-read proteogenomics

Madison M. Mehlferber^{a,b}, Erin D. Jeffery^b, Jamie Saquing^b, Ben T. Jordan^b, Leon Sheynkman^b, Mayank Murali^b, Gael Genet^c, Bipul R. Acharya^{c,d,e}, Karen K. Hirschi^{c,d}, and Gloria M. Sheynkman ^{a,b,f,g}

^aDepartment of Biochemistry and Molecular Genetics, University of Virginia, Charlottesville, VA, USA; ^bDepartment of Molecular Physiology and Biological Physics, University of Virginia, Charlottesville, Virginia, USA; ^cDepartment of Cell Biology, University of Virginia School of Medicine, Charlottesville, VA, USA; ^dCardiovascular Research Center, University of Virginia, Charlottesville, VA, USA; ^eWellcome Centre for Cell-Matrix Research, Faculty of Biology, Medicine and Health, the University of Manchester, UK; ^fCenter for Public Health Genomics, University of Virginia, Charlottesville, VA, USA; ^gUVA Comprehensive Cancer Center, University of Virginia, Charlottesville, Virginia, USA

ABSTRACT

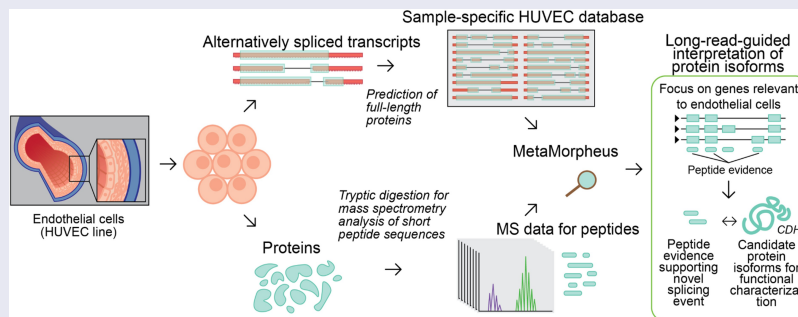
Endothelial cells (ECs) comprise the luminal lining of all blood vessels and are critical for the functioning of the cardiovascular system. Their phenotypes can be modulated by alternative splicing of RNA to produce distinct protein isoforms. To characterize the RNA and protein isoform landscape within ECs, we applied a long read proteogenomics approach to analyse human umbilical vein endothelial cells (HUVECs). Transcripts delineated from PacBio sequencing serve as the basis for a sample-specific protein database used for downstream mass-spectrometry (MS) analysis to infer protein isoform expression. We detected 53,863 transcript isoforms from 10,426 genes, with 22,195 of those transcripts being novel. Furthermore, the predominant isoform in HUVECs does not correspond with the accepted “reference isoform” 25% of the time, with vascular pathway-related genes among this group. We found 2,597 protein isoforms supported through unique peptides, with an additional 2,280 isoforms nominated upon incorporation of long-read transcript evidence. We characterized a novel alternative acceptor for endothelial-related gene *CDH5*, suggesting potential changes in its associated signalling pathways. Finally, we identified novel protein isoforms arising from a diversity of RNA splicing mechanisms supported by uniquely mapped novel peptides. Our results represent a high-resolution atlas of known and novel isoforms of potential relevance to endothelial phenotypes and function.

ARTICLE HISTORY

Received 4 August 2022
Revised 20 October 2022
Accepted 26 October 2022

KEYWORDS




Long-read RNA-seq; PacBio; isoforms; mass-spectrometry-based proteomics; proteogenomics; cardiovascular; endothelial cells; HUVECs; alternative splicing; protein isoforms; nextflow; MetaMorpheus; Orbitrap



Introduction

Endothelial cells are critical for the development and maintenance of the cardiovascular system. They form the lining of all blood vessels within the body allowing for functions such as oxygen nutrient delivery, blood pressure regulation, and immune control [1]. Endothelial dysfunctions can contribute to a host of cardiovascular diseases, such as atherosclerosis, diabetes retinopathy, cancer, and stroke [2]. Improved understanding of these and related diseases may be attained through molecular characterization of the proteome underlying endothelial cell identity and functionality [3,4].

Endothelial cells can express functionally distinct protein isoforms through the process of alternative splicing (AS). For example, vascular endothelial growth factor A (VEGF-A) exists as two separate isoform families that differentially bind to the extracellular region on VEGFR1 or VEGFR2 leading to proliferation and survival of endothelial cells. One VEGF-A isoform family is pro-angiogenic and another is anti-angiogenic [5,6]. Together these isoforms work in balance to regulate new vessel formation. Globally, across the endothelial cell proteome, many gene functions are modulated by AS [7–11]. However, despite many high-throughput sequencing datasets collected on endothelial cells [12],

CONTACT Gloria M. Sheynkman  gs9yr@virginia.edu  The Department of Molecular Physiology and Biological Physics, University of Virginia, Charlottesville, VA, USA
 Supplemental data for this article can be accessed online at <https://doi.org/10.1080/15476286.2022.2141938>

© 2022 The Author(s). Published by Informa UK Limited, trading as Taylor & Francis Group.
This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted use, distribution, and reproduction in any medium, provided the original work is properly cited.

our knowledge of individual protein isoforms that are expressed is incomplete [13].

In order to characterize the proteome of endothelial cells, human umbilical vein endothelial cells (HUVECs) can serve as a relevant model system, since they are primary cells that can be expanded in culture to generate sufficient material for proteomic analysis [3,14,15]. A prior study performed by Madugundu and colleagues employed a proteogenomics approach, incorporating RNA-seq and mass-spectrometry (MS)-based proteomics in order to characterize proteomic variation in HUVECs [16]. By utilizing short-read RNA-seq data, the authors generated a set of custom databases of relevance to protein variants. Though the main focus of the study was to characterize diverse sources of variation, such as single amino acid variants and phosphorylations, they generated a database of candidate splice-junction peptides derived from novel exon-to-exon connections (i.e., junctions), as well as a custom database based on inferred reconstruction of full-length transcripts. The study reported a few novel splice junction peptides, providing further insight into the role of splicing events in HUVECs. However, the proteogenomics approach used relied upon short-read RNA sequencing in the custom database generation, and short reads cannot provide unambiguous knowledge of the bona fide full-length isoform (i.e. complete chain of exon/junction connectivities) [17], which is needed for accurate prediction and detection of full-length protein isoforms [18].

For improved characterization of protein isoform expression in HUVECs, it would be ideal to obtain full-length transcript information to infer expressed isoforms at the protein level. Fortunately, advances in sequencing technology, such as through the PacBio or Oxford Nanopore long-read sequencing platforms, have allowed for detection of full-length transcript isoforms [19–21]. Capitalizing on these technologies, we previously developed a proteogenomic approach that incorporates long-read RNA sequencing with MS analysis, which we term ‘long-read proteogenomics’ [22]. Long-read RNA-seq returns information on full-length transcript isoforms [23], which is bioinformatically translated into full-length protein isoform predictions [22,24–26]. These predicted protein isoforms serve as sample-specific, full-length isoform models from which to infer protein expression from MS data [27].

Here, we apply a long-read proteogenomic approach to characterize protein isoforms expressed in HUVECs. We demonstrate the application of PacBio long-read RNA-seq data towards characterization of the full-length transcriptome in HUVECs, which includes detection of unannotated transcript isoforms. A PacBio-derived HUVEC protein database is searched against a sample-matched MS dataset facilitating the characterization of HUVEC-specific isoforms. Finally, we report on the discovery of novel peptides, providing evidence for novel isoforms through a direct mapping of novel peptides to full-length protein isoforms in HUVECs. Overall, we present the first application of a long-read proteogenomics approach as applied to primary endothelial cells. These results nominate candidate isoforms for functional studies of how splicing modulates endothelial cell phenotype and function.

Experimental methods

HUVEC cell culture

Primary Human Umbilical Vein Endothelial Cells (HUVECs) were purchased from Lonza (C2519AS) and used up to passage five. Early passage HUVECs were cultured in EGM™2-Bulletkit™ medium with growth supplements CC-3156 & CC-4176 purchased from Lonza. At 80% confluency, HUVECs were trypsinized, washed twice with phosphate-buffered saline (PBS), pelleted, and frozen at –80°C.

Long-read RNA-seq (PacBio Iso-Seq) library preparation and sequencing run

PacBio (Iso-Seq) data were collected on the extracted total RNA collected from the HUVEC cell pellet. HUVEC RNA was analysed on an Agilent Bioanalyzer to confirm concentration and RNA integrity for downstream analysis. We observed a RIN value of 10. From this RNA, cDNA was synthesized using the NEB Single Cell/Low Input cDNA Synthesis and Amplification Module (New England Biolabs).

Approximately 200 ng of HUVEC cDNA was converted into a SMRTbell library for usage with the Iso-Seq Express Kit SMRT Bell Express Template prep kit 2.0 (Pacific Biosciences). Through this protocol, bead-based size selection occurs in order to remove low mass cDNA (less than 500 kb). Each SMRTbell library was sequenced on the SMRT cell on Sequel II system. A 2-hour extension and 3-hour movie collection time was used for data collection. The ‘ccs’ command from the PacBio SMRTLink suite (SMRTLink version 9) was used to convert raw reads into Circular Consensus (CCS) reads.

Mass spectrometry-based proteomics sample preparation

Harvested HUVECs, approximately 5 million cells each, were pelleted and frozen at –80°C. The sample pellet was lysed according to the Filter Aided Sample Preparation (FASP) protocol [28]. Lysis buffer used in the FASP was changed to 6% SDS, 150 mM DTT, 75 mM Tris-HCl. To the 30 µL pellet of 5 million cells, an aliquot of 60 µL of lysis buffer was added and probe-sonicated to lyse the cells and shear the nucleotide material. Sonication continued for 1–5 minutes until the sample was clear and no longer viscous. The lysate was then incubated at 95°C for 5 minutes. Protein quantitation was estimated by BCA assay to be approximately 500–600 µg. Quadruplicate aliquots of 20 µL each were subjected to FASP and trypsin digest (1 µg per aliquot) and allowed to incubate at 37°C overnight. Nanodrop analysis estimated peptide content at 22 µg per trypsin digest (total of 88 µg).

Offline HPLC Fractionation

The tryptic digests were pooled and dried down to a volume of 40 µL and subjected to offline high pH RP-HPLC fractionation using an Agilent 1200 HPLC. Sample was loaded onto a Thermo Scientific Hypersil Gold C18 column (150 mm × 3 mm × 3 µm C18), equilibrated with 95% solvent A (20 mM NH₄ formate, pH 10) and 5% solvent B (70% acetonitrile/30%

solvent A), and eluted at a flow rate of 400 $\mu\text{L}/\text{min}$, with fractions collected every 1 minute from RT 38–63 min. The following gradient was used: 5% B from 0 to 30 min, 5–65% B from 30 to 63 min, 65–100% B from 64 to 69 min, 100–5% B from 69 to 70 min, 5% B from 70 to 73 min. Samples containing peptide, according to UV 214 nm corresponding to the HUVEC pellet were digested with trypsin. Collected fractions 4–20 were selected for LC-MS/MS analysis.

NanoLC-MS/MS analysis

The resulting peptides were dried to 12 μL and analysed by nanoLC-MS/MS using a Dionex Ultimate 3000 (Thermo Fisher Scientific, Bremen, Germany) coupled to an Orbitrap Eclipse Tribrid mass spectrometer (Thermo Fisher Scientific, Bremen, Germany). Three microlitres of each peptide-containing sample were loaded onto an Acclaim PepMap 100 trap column (300 $\mu\text{m} \times 5 \text{ mm} \times 5 \mu\text{m}$ C18) and gradient-eluted from an Acclaim PepMap 100 analytical column (75 $\mu\text{m} \times 25 \text{ cm}$, 3 μm C18) equilibrated in 96% solvent A (0.1% formic acid in water) and 4% solvent B (80% acetonitrile in 0.1% formic acid). The peptides were eluted at 300 nL/min using the following gradient: 4% B from 0 to 5 min, 4–28% B from 5 to 210 min, 28–40% B from 210 to 240 min, 40–95% B from 240 to 250 min and 95% B from 250 to 260 min.

The Orbitrap Eclipse was operated in positive ion mode with 1.9 kV at the spray source, RF lens at 30% and data dependent MS/MS acquisition with XCalibur version 4.3.73.11. Positive ion Full MS scans were acquired in the Orbitrap from 375 to 1500 m/z with 120,000 resolution. Data dependent selection of precursor ions was performed in Cycle Time mode, with 3 seconds in between Master Scans, using an intensity threshold of 2×10^4 ion counts and applying dynamic exclusion ($n = 1$ scans within 30 seconds for an exclusion duration of 60 seconds and ± 10 ppm mass tolerance). Monoisotopic peak determination was applied and charge states 2–6 were included for HCD scans (quadrupole isolation mode; 1.6 m/z isolation window). The resulting fragments were detected in the Orbitrap at 15,000 resolution with standard AGC target.

Long-read RNA-seq analysis, MS searching, and proteogenomic analysis conducted using a Nextflow pipeline

The long-read proteogenomics pipeline was implemented with Nextflow, a workflow framework which allows for scalable and reproducible computational analysis. The Nextflow pipeline developed and described previously was used to process HUVEC collected PacBio data, translate the resulting transcripts into the protein database (see Deriving a HUVEC sample-specific protein isoform database below), and perform proteomics database searches [22]. Further information on the workflow including individual modules of the Nextflow pipeline can be found at <https://github.com/sheynkman-lab/Long-Read-Proteogenomics> [22]. The GitHub revision (i.e., commit) used in this analysis was <https://github.com/sheynkman-lab/Long-Read-Proteogenomics/releases/tag/v1.0.0>. All

transcriptomic and proteogenomic docker images that are used within the analysis can be found at <https://hub.docker.com/> under the repository `gsheynkmanlab`. The analysis was performed on the University of Virginia High Performance Computing system.

Long-read RNA-seq (PacBio Iso-Seq) data analysis

The CCS reads obtained from PacBio sequencing were analysed using the IsoSeq workflow described previously [22]. Primer removal was done on the 5' and 3' end. The 5' primer consists of an NEB adapter sequence (Sequence: GCAATGAAGTCGCAGGGTTGGG). The 3' primer consists of the Clontech SMARTer cDNA universal primer (Sequence: GTACTCTGCGTTGATACCACTGCTT). Following processing of the raw reads using the IsoSeq workflow, we derived the number of full-length reads corresponding to each distinct transcript. Full-length read counts per million (CPM) were computed by dividing the number of full-length reads aligning to a transcript isoform by the total number of reads and then multiplying this by a factor of 1,000,000.

Transcript isoform classification and filtering

SQANTI is a computational tool used for comparison, classification, and quality assessment of the full-length isoforms sequences collected from the long-read platform [29]. We used SQANTI3 (version 1.3) to annotate the polished transcript isoforms obtained from the Iso-Seq analysis using SQANTI default parameters. Note: the default parameters included options to use the genome-derived sequences for the isoform output. As a result, transcriptional variations inclusive of alternative N-termini, alternative splicing, etc. but not genetic variations are captured in the HUVEC sample-specific database.

Generation of a full-length protein isoform database from the long-read RNA-seq data

After deriving a high confidence set of full-length transcript isoforms, within the Nextflow pipeline we select the most biologically plausible ORF for each of the Iso-Seq transcripts. Calling the best ORF consists of two steps: finding candidate ORFs (50 nucleotides or longer) using CPAT [30], and selecting the most plausible ORF based on coding potential, relation of AUG start site to GENCODE reference start sites, and number of AUGs skipped to reach the ORF start site.

To generate the PacBio-derived protein database (HUVEC sample-specific database) employed for downstream MS searching, transcripts were grouped that produced ORFs of the same sequence. The total transcript abundance for each grouping was calculated as the sum of all CPM values for the transcripts comprising that group. Candidate isoforms are further classified based on the protein sequence in relation to the reference protein isoforms, as defined in the 'sqanti_protein' and 'protein_classification' modules in the Nextflow pipeline. Classifications are based on a variant of nomenclature used within the SQANTI3 software, which we call 'SQANTI Protein'.

Additional filtering was performed in order to retain only isoforms that were likely protein coding. Isoforms that did not have a stop codon within the predicted ORF, and could represent truncations, were removed. Isoforms that were either fully mapped to a protein-coding GENCODE reference isoform ('protein full splice match', pFSM) were retained, as well as isoforms that contained a novel combination of known splice sites or junctions (pNIC). Of the isoforms that contain novel splice sites (pNNC), suspected nonsense mediated decay (NMD) isoforms were removed. Here, NMD suspects were defined as isoforms that contained more than two junctions after the stop codon. Isoforms that were not classified as pFSM, pNIC, or pNNC were removed from consideration. Protein classification details can be found within the 'protein_classification' module of the pipeline, while the filtering criteria can be found within the 'protein_filter' module of the Nextflow pipeline.

A hybrid database was developed that incorporated isoforms from PacBio if the gene resided in the high confidence region, defined as where the aggregated transcriptomic gene abundance contained at least three CPM and the average reference transcript length was between 1 and 4 kilobases (kbp). If a gene did not meet these criteria, the reference isoforms were substituted in place of the long-read isoforms. If a gene was not found within the long-read transcriptomic data, the reference protein isoforms were also appended into the hybrid database. A detailed description of reasoning behind creation of a hybrid database has been described previously [22].

GENCODE and UniProt reference protein database

The GENCODE protein database used in this study was created by downloading the coding translation FASTA and grouping entries with the same protein sequence for each gene ('make_gencode_database' module in the Nextflow). For the many cases where one or more GENCODE transcripts from the same gene lead to the same protein sequence, the transcripts were grouped and assigned a protein accession as the first alphanumeric GENCODE protein accession, by the transcript name (e.g., GAPDH-201).

The UniProt database used was the reviewed human database with isoforms, downloaded November 1st, 2020. The database contains 42,358 protein isoform entries from 20,292 genes.

MS database search

Standard proteomic analysis of acquired mass spectra files were performed using the free and open source search software program MetaMorpheus [31]. A custom branch and Docker image were made as part of the Nextflow pipeline (GitHub: <https://github.com/smith-chem-wisc/MetaMorpheus/tree/LongReadProteogenomics>, Docker: https://hub.docker.com/r/smithchemwisc/metamorpheus/tags?page=1&ordering=last_updated tag; lrproteogenomics) based on MetaMorpheus version 0.0.316. Analysis of the collected spectra files performed either using the HUVEC sample-specific database (HUVEC-derived PacBio reads + GENCODE entries; 'HUVEC sample-specific

database') (71,511 of entries from 19,982 genes) in which the subset of PacBio derived entries are 26,675 protein isoforms from 7,283 genes. The GENCODE human database (version 35; 87,729 protein entries from 19,982 genes), or the UniProt reviewed human database with isoforms (downloaded 8 July 2021; 42,380 protein entries from 20,292 genes). All searches were conducted with a contaminants database, included in MetaMorpheus, which contains 264 common contaminant proteins frequently found in MS samples.

All RAW spectra files were first converted to mzML format with MSConvert prior to analysis with MetaMorpheus (see 'mass_spec_raw_convert' module in the Nextflow pipeline). For the MetaMorpheus MS search, the settings used for all search tasks can be found in Supplementary Information Table S6. MetaMorpheus produces peptide spectral match (PSM), peptide and protein group result files, which we analysed in downstream custom modules. All peptide and protein results reported employ a 1% False Discovery Rate (FDR) threshold after target-decoy searching [32].

Criteria for Novel Peptide Identification

Stringent filtering criteria and manual validation were used, as described previously [22,33] to ensure that the spectrum does in fact represent the novel peptide sequence. Spectra corresponding to the scan number of the identified novel peptide sequence were derived from MetaDraw and manually inserted into an Excel file which were then manually evaluated. Corresponding University of California Santa Cruz Genome Browser tracks depicting protein isoforms were derived and can be found via the following session: https://genome.ucsc.edu/s/mm5db/211018_huvec_hcd_trp. In addition to previously [22] described criteria for novel peptide annotation, we allowed for cases where the C13 isotope for a novel peptide was selected as the precursor.

Data analysis and plot generation

All downstream data analyses were performed through custom Python scripts. Data analysis scripts used for generation of figures, plots, and statistics may be found in the following GitHub repository: <https://github.com/sheynkman-lab/Huvec-Proteogenomic-Analysis>

Availability of data and materials

Raw long-read RNA-seq data collected on the PacBio platform are available from the Sequence Read Archive (PRJNA832812, corresponding to accession SRR18959149). Data generated by mass spectrometry are available through MassIVE, the Mass Spectrometry Interactive Virtual Environment (MSV000089326). The output of the data analysis including the long-read proteogenomics Nextflow workflow results generated using the mass spectrometry and long-read RNA-sequencing data as well as the post pipeline analysis results are available on Zenodo (<https://zenodo.org/record/7117445#.Y2FQE-wpD0o>).

The open-source software produced in the making of this work is freely available under the MIT licence found in the GitHub repository (<https://github.com/sheynkman-lab/Long-Read-Proteogenomics>). A wiki was created (<https://github.com>).

[com/sheynkman-lab/Long-Read-Proteogenomics/wiki](https://github.com/sheynkman-lab/Long-Read-Proteogenomics/wiki)) describing each of the pipeline processes.

Code used to generate the main figures and tables in this manuscript can be found in the GitHub repository (<https://github.com/sheynkman-lab/Huvec-Proteogenomic-Analysis>).

Results

Long-read proteogenomics to characterize isoforms in endothelial cells

In order to characterize the isoforms expressed in an endothelial cell population, we subjected HUVECs to ‘long-read proteogenomics’ where samples undergo long-read RNA-sequencing and mass-spectrometry analysis in parallel, which is followed by integrative analysis of the matched datasets [22]. The full-length transcripts – obtained from long-read RNA-seq – are converted to a predicted protein database, serving as candidate isoforms for proteomic detection (Fig. 1). As a first step in our method, PacBio RNA sequencing is performed to characterize the HUVEC transcriptome.

Long-read RNA-seq of HUVECs reveals widespread and novel isoform diversity

Long-read RNA-seq data was collected on the PacBio sequencing platform using the ‘Iso-Seq’ method [34], generating 3,608,972 long-reads (i.e. circular consensus reads). These reads were processed by Iso-Seq3 [34] to generate the set of distinct transcript isoforms and their respective abundances (Fig. 2A)[22].

PacBio-derived transcripts were compared to reference transcripts (GENCODE v35) and their novelty status was defined using SQANTI3 (Fig. 2B)[29]. The UniProt database lacks a complete mapping of protein isoforms to the reference genome, and therefore we could not compare transcripts to UniProt directly, although future efforts may address this limitation [35–37]. Based on a comparison to GENCODE models, we identified 53,863 transcripts from 10,426 protein coding genes, inclusive of all transcripts with a minimum

abundance of one full-length read count per million (CPM). The average length of transcripts is 2,846 kilobase pairs (kbp) (Supplementary Information Figure S1A). Among the 53,863 transcripts isoforms identified in the HUVEC sample, 31,668 (59%) matched exactly to a transcript isoform in GENCODE, the match being based on splice junction connectivity (‘full splice matches’, ‘FSM’, Fig. 2B). The remaining 22,195 (41%) isoforms were unannotated, or novel, in terms of the observed ordering of splice junctions along the length of the transcript (Fig. 2B). Of the unannotated isoforms identified, 13,746 (62%) contained novel combinations of known splice junctions (‘novel in catalog’, ‘NIC’), and the remaining 8,449 (38%) isoforms contained entirely new exon splice boundaries, in which the acceptor or donor site is not represented in GENCODE (‘novel not in catalog’, ‘NNC’, Fig. 2B). The overall abundance distribution for identified transcripts was wide ranging (see Supporting Information Figure S1B). As expected, on average, the novel transcripts exhibit lower abundance than known transcripts (Fig. 2C) [38,39]. The FSM transcripts displayed a median abundance of 2.4 CPM, while the NIC and NNC transcripts displayed a median of 1.5 and 1.3 CPM, respectively. These data illustrates that novel transcripts tend to exhibit lower abundances than known transcripts. While these trends represent average expression differences, particular novel transcripts can exhibit high abundances within HUVECs.

Using the full-length transcriptomics dataset, we next determined the number of protein-coding genes that returned evidence for expression of multiple isoforms. We found that 82% (8,522 genes of the 10,426 genes represented) of detected genes expressed multiple transcript isoforms (Fig. 2D). To focus on genes involved in endothelial pathways that may be co-expressing multiple isoforms, we manually curated the literature to compile a list of genes that are involved in vascular pathways related to early endothelial differentiation and development or hemogenic specification (see Supporting Information Table S1) [40,41]. We then determined which endothelial genes are expressing multiple isoforms in our HUVEC sample. To have increased confidence in isoform expression of such genes, we filtered for genes which contain

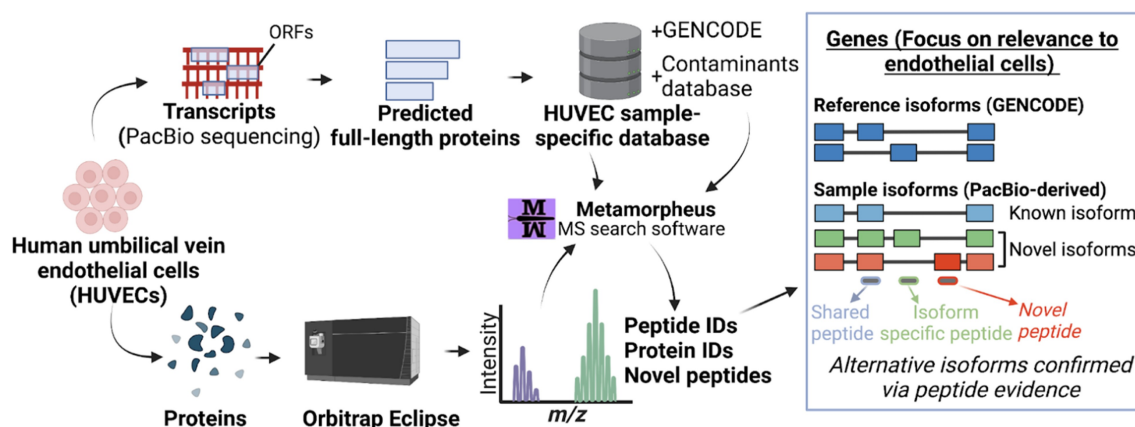


Figure 1. Characterization of isoform diversity in HUVECs through integration of long-read RNA-seq with mass spectrometry data (‘long read proteogenomics’). Transcripts are converted into a protein isoform database based on predicted open reading frames (ORFs) and the resulting database is searched against a sample-matched bottom-up mass spectrometry (MS) dataset. The peptide identifications can be used to support the expression of isoform candidates related to endothelial pathways.

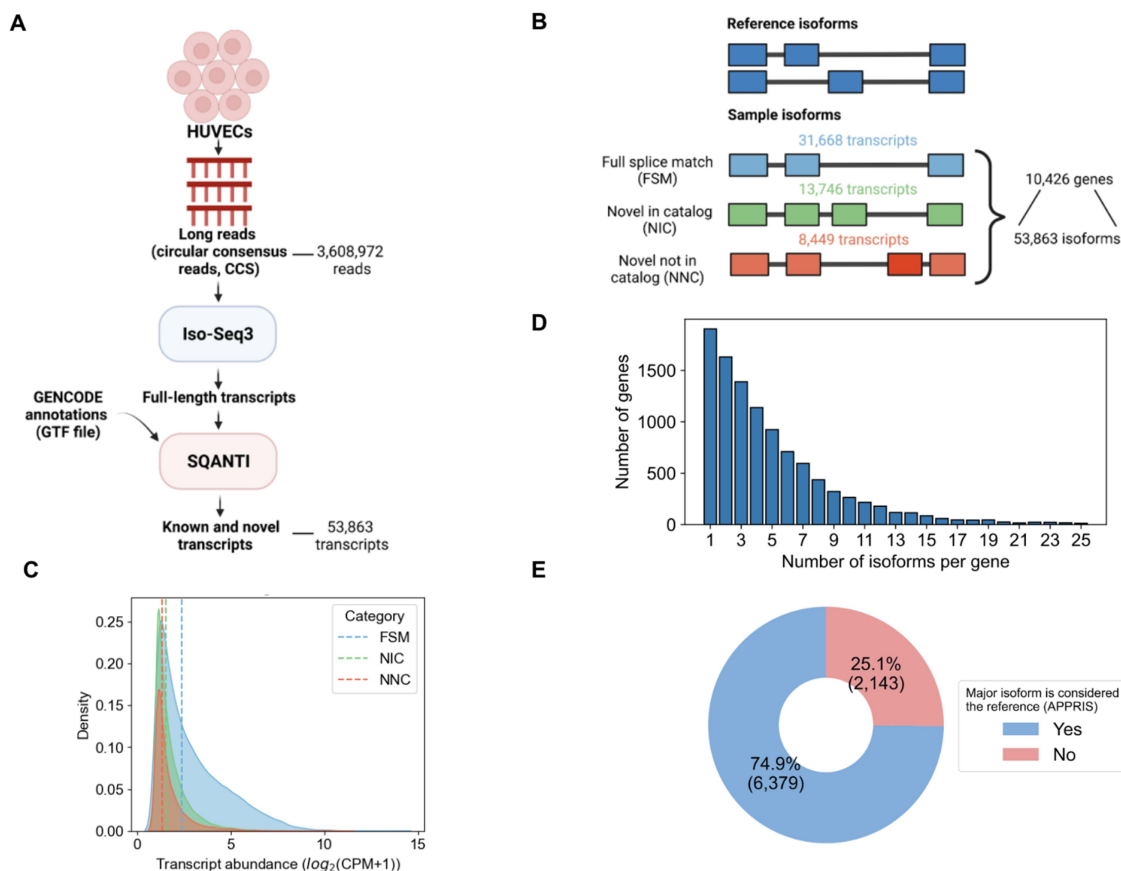


Figure 2. Characterization of transcript isoform diversity in HUVECs via long-read RNA-Seq. (A) Schematic of the long-read RNA-seq analysis pipeline. (B) Transcripts and genes identified from PacBio long-read RNA-seq. The number of known (blue) and novel isoforms (green and orange) are shown. (C) Transcript abundance distribution for known (FSM) versus novel transcripts (NIC, NNC), with dashed lines representing median abundance values in full-length read counts per million (CPM) for each category (FSM = 2.4, NIC = 1.5, NNC = 1.3). (D) Distribution of the number of genes expressing multiple isoforms. (E) Fraction of genes in which the most abundantly expressed isoform (“major isoform”) differs from the reference isoform (APPRIS principal isoform).

two or more isoforms with each isoform having an abundance of at least three CPM. We identified multiple co-expressing isoforms for *CD34*, *CELFI*, *FLT1*, *NRP1* and *SRSF5* (Table 1, with annotations from GOrilla [42]) [6,43].

To explore the putative functional effects of candidate genes, we closely examined the potential impacts of changes to amino acid sequences among isoforms of *NRP1*, *CELFI* and *FLT1*. We discovered novel isoforms for *NRP1*, also called neuropilin. *NRP1* is involved in regulating angiogenesis and arteriogenesis pathways through its binding interactions with VEGF-A [40,43,44]. Notably, we detected a novel isoform (PB.6952.10) at moderate abundance (40 CPM) containing an alternative donor region. This region has been found within three amino acids of a glycosylation site that has been suggested as potentially affecting neuropilin activity [6]. Additionally, we found *NRP1* isoform expression of both soluble and membrane-bound forms, and it has been well known that the soluble form acts antagonistically to the full-length form for VEGF signalling [6]. Finally, we identified an isoform with a skipped exon in the C-terminal disordered region of the protein that resides just outside of the trans-membrane domain.

Next, we examined *CELFI*, which is an RNA binding protein that is a known regulator of splicing in cardiovascular biology [45]. We observed eleven isoforms for this gene,

seven being novel. The abundance for the major isoform of *CELFI* is moderately high (124 CPM) (PB.7605.2), but the 2nd to 5th isoforms by ranked abundance are also expressed at moderate levels, with two of them novel (PB.7605.5, PB.7605.1). These novel *CELFI* isoforms arise from distinct combinations of splicing events. Nearly all isoforms contain the complete set of three RNA recognition motif (RRM) domains, as described previously [46]; however, an alternative acceptor site residing between the 2nd and 3rd RRM domain may alter the inter-domain distance, which may alter binding behaviour. Interestingly, the *CELFI* isoforms contain either an extended or truncated N-terminus, which may have a direct effect on cellular localization – based on previous reports of the extended N-terminal *CELFI* isoform as being localized to the nucleus and the truncated N-terminal *CELFI* isoform as being localized to the cytoplasm [47]. Based on the long-read RNA-seq data, we estimate that in HUVECs, approximately 30% of *CELFI* isoforms may be localized to the nucleus.

Finally, we examined *FLT1*, a gene that encodes the vascular endothelial growth factor receptor 1 (VEGFR1) and mediates VEGF-A signalling allowing for the survival and proliferation of endothelial cells [40]. We identified ten protein isoforms for *FLT1*. Five of such isoforms were novel and were all extremely low in abundance (~1 CPM, only a few

Table 1. Endothelial-relevant genes expressing multiple transcript isoforms in HUVECs.

Gene	PacBio transcript	GENCODE isoform match	Counts per million (CPM)	Function*
<i>CD34</i>	PB.1222.12	CD34-201	13.9	Cell adhesion molecule
	PB.1222.24	novel	9.8	
	PB.1222.26	CD34-202	226.7	
	PB.1222.27	CD34-201	344.7	
	PB.1222.29	novel	12.8	
<i>CELF1</i>	PB.1222.31	CD34-203	27.8	Pre-mRNA splicing
	PB.7605.14	novel	9.4	
	PB.7605.21	novel	4.5	
	PB.7605.23	novel	6.8	
	PB.7605.28	novel	58.3	
	PB.7605.29	CELF1-201	13.2	
	PB.7605.56	novel	5.3	
	PB.7605.70	novel	27.1	
	PB.7605.76	novel	47.0	
	PB.7605.81	novel	10.1	
<i>CDH5</i>	PB.10443.1	CDH5-209	216.5	Regulation of cellular metabolic process
	PB.10443.11	CDH5-201	3.0	
	PB.10443.15	novel	321.8	
	PB.10443.18	novel	3.0	
	PB.10443.2	CDH5-201	2402.9	
	PB.10443.22	CDH5-201	44.0	
	PB.10443.26	novel	9.0	
	PB.10443.28	CDH5-201	13.5	
	PB.10443.33	novel	4.5	
	PB.10443.36	CDH5-201	20.7	
	PB.10443.38	CDH5-201	28.2	
	PB.10443.40	CDH5-208	5.6	
	PB.10443.45	novel	27.1	
	PB.10443.49	novel	9.8	
	PB.10443.50	CDH5-209	6.4	
	PB.10443.52	CDH5-201	10.5	
	PB.10443.57	novel	3.0	
<i>FLT1</i>	PB.8882.15	FLT1-204	42.9	Vascular endothelial growth factor activated receptor activity
	PB.8882.22	FLT1-207	23.7	
	PB.8882.27	FLT1-207	9.4	
	PB.8882.30	FLT1-207	29.7	
	PB.8882.9	FLT1-201	6.4	
<i>NRP1</i>	PB.6952.10	novel	32.0	Vascular endothelial growth factor binding
	PB.6952.12	novel	3.8	
	PB.6952.35	novel	4.1	
	PB.6952.54	novel	3.4	
	PB.6952.58	novel	5.6	
<i>PECAM1</i>	PB.11293.22	novel	29.3	Epithelium development
	PB.11293.23	novel	38.7	
	PB.11293.54	novel	6.0	
	PB.11293.55	novel	5.3	
	PB.11293.64	PECAM1-203	524.8	
	PB.11293.68	novel	32.7	
	PB.11293.7	novel	6.4	
	PB.11293.70	novel	12.8	
	PB.11293.71	novel	4.5	
	PB.11293.80	novel	5.6	
	PB.11293.81	novel	29.3	
	PB.11293.83	novel	3.4	
	PB.11293.9	novel	5.3	
	PB.11293.95	novel	3.8	
	PB.11293.98	novel	3.4	
<i>SRSF5</i>	PB.9356.16	SRSF5-201	10.5	Pre-mRNA splicing
	PB.9356.17	SRSF5-217	103.8	
	PB.9356.21	SRSF5-217	68.4	
	PB.9356.4	SRSF5-207	15.8	

***Function** – GO annotations derived from GOrilla

reads supporting their existence); therefore, we did not consider them further. Among all the isoforms, we observed two major families of *FLT1* isoforms: 1) full-length isoforms that contain the transmembrane domain, and can promote endothelial proliferation and angiogenesis [9], and 2) short, soluble isoforms that lack the transmembrane domain but still binds to VEGF-A, and thus loses its signal transduction function, and therefore is anti-angiogenic [6].

Given the prevalence of genes that co-express multiple isoforms in HUVECs, we next asked to what extent the

identity of the most highly expressed isoform, i.e. the major isoform, match what is defined as the ‘reference isoform’ for a gene. To define a gene’s ‘reference isoform’, we used the APPRIS database which reports a principal isoform to be most representative for a gene [48]. The APPRIS principal isoform concept is related to the concept of a UniProt ‘canonical’ protein, though underlying assumptions differ [36,48]. For the genes expressing multiple isoforms, we classified their corresponding isoforms as either major, i.e. the most abundant isoform based on relative expression

levels of all isoforms for a gene, or minor. There were 1,904 genes only expressing one isoform and therefore were excluded from analysis. We identified 8,522 transcripts as the major isoform and 43,437 as minor isoforms. We found, as expected, that on average the major isoforms are more highly expressed than minor isoforms (Supplementary Information Figure S1C-D). Surprisingly, we found that for 25% (2,143 isoforms) of genes, the major isoforms in our HUVEC sample do not coincide with the APPRIS principal isoform (Fig. 2E). Within this population of major isoforms, we found six genes involved in endothelial pathways, *CELF1*, *FLT1*, *GATA2*, *NR2F2*, *NRP1*, *NRP2* and *SRSF6* (see Supporting Information Table S2). These results illustrate that the major isoform expressed in a given sample may not always correspond to the generic “reference” isoform for a gene, which can be explained by the fact that isoforms exhibit cell or tissue-specific expression patterns [49].

Next, we examined the presence of previously annotated splice factors [50] expressed within our HUVEC PacBio data. Overall, we detected long reads for 85 annotated splice factors, with the 10 most abundant splice factors including *HNRNPA2B1*, *HNRNPK*, *HNRNPC*, *DDX5*, *EWSR1*, *PCBP2*, *HNRNPA1*, *PCBP1*, *FUS*, *KHDRBS1* (Supplementary Table S3). Notably, *SRSF5* was found as the eleventh highest expressed splice factor at 408 CPM, followed by *SRSF2* as the twenty-second most abundant splice factor at 300 CPM and lastly *CELF1* as the 25th highest expressed at 263 CPMs. Interestingly, it has been observed that *SRSF2* and *SRSF5* are involved in splicing of VEGF-A pre-mRNA splicing [5].

We next asked how the novel isoforms differed from the APPRIS principal isoform in terms of length and affected amino acids. As expected, on average, novel isoforms are shorter than the reference form due to the loss of amino acid regions (Supporting Information Figure S2C-D), with a median shortening of 159 amino acids and an average gain of 11 amino acids. The APPRIS principal isoform for a gene may not be the most representative isoform in HUVECs (see Fig. 2E). Therefore, we also compared the lengths of the novel isoform against the ‘major’ isoform in HUVECs, i.e. the highest expressed isoform in the HUVEC data. Interestingly, we observed that ‘major’ isoforms do not tend to be the longest isoform of a gene, or at least this trend is not as stark as with APPRIS principal isoforms. This is likely because the APPRIS algorithm does have a tendency to select the longest isoform of a gene as its ‘principal’ isoform [48].

Collectively, the HUVEC transcriptomic results demonstrate the use of long-read RNA-seq to characterize sample-specific variation in isoform identity and abundance.

Deriving a HUVEC sample-specific protein isoform database

The vast transcriptome diversity of HUVECs likely translates in some part to a diversity of protein isoforms. To explore this question, we translated the HUVEC transcript isoform sequences *in silico* into open-reading frames (ORFs) and compiled the predicted sequences into a HUVEC sample-specific

protein isoform database for MS searching, as previously described (Supplementary Information Figure S2A) [22]. To classify the relationships between the predicted proteins to that of annotated protein isoforms in GENCODE [35], we used the classification scheme we previously developed, SQANTI Protein [22]. SQANTI Protein automatically categorizes known and novel protein isoforms. The categories include ‘protein full-splice match’ (pFSM), ‘protein novel in catalog’ (pNIC), and ‘protein novel not in catalog’ (pNNC) (Supplementary Information Figure S2B). We found that 16,296 predicted proteins exactly matched protein isoforms in the GENCODE reference (pFSMs), while 24,896 predicted protein isoforms were novel (Supplementary Information Figure S2C). Among those novel isoforms, 5,855 had novel combinations of known protein sequence elements such as the N-terminus, the C-terminus or the splicing pattern (pNICs). The other 19,041 protein isoforms had one or more entirely novel elements, such as a novel N-terminus or an unannotated exon (pNNC).

Among the candidate protein isoforms, we first filtered out protein isoforms that may have resulted from transcripts from incomplete reads or poor-quality transcripts (see *Protein database generation* in Methods; 11,876 filtered out). The remaining 34,531 predicted protein isoforms (comprising 16,296 pFSMs, 5,855 pNICs, and 12,389 pNNCs) from 10,912 genes were compiled to create a preliminary HUVEC protein database (Fig. 3A). These genes and their associated isoforms represent candidates for inclusion in the final database. For the final database, we decided to only include isoforms from genes for which we could ensure a complete sampling of the transcripts, and thus the predicted proteins. Therefore, we created a hybrid database in which we defined a core set of genes for which the transcript detection, and thus predicted proteins, is likely complete based on the long-read data collected. The core set of genes included in the hybrid database have a minimum abundance of three CPM and a moderate transcript length (1–4 kbp average GENCODE-annotated transcript length). For all other genes, the hybrid database is populated with all GENCODE protein isoform entries. The hybrid structure of the final database ensures comprehensiveness of the protein models, with the protein completeness assumption of target-decoy searching satisfied so as to avoid issues of an off-target peptide match [32].

As described, the final HUVEC sample-specific database for proteomic analysis includes a mixture of custom PacBio-derived proteins as well as annotated GENCODE proteins (Table 2). A detailed listing of steps to convert the transcriptome data to a protein database may be found in Supplementary Information Table S4.

Collection of a deep-coverage MS dataset for HUVECs

In order to characterize protein isoforms in HUVECs, we generated and analysed a deep-coverage MS dataset collected on the same HUVEC pellets that were used for long-read RNA sequencing (Fig. 3B). HUVECs were lysed and processed using the filter aided sample preparation (FASP) protocol, in which protein was digested with trypsin to generate a mixture of tryptic peptides. The tryptic digest was subjected to off-line fractionation on an analytical scale high-pH reverse-phase liquid chromatography instrument, and 20 fractions were collected (Supplementary

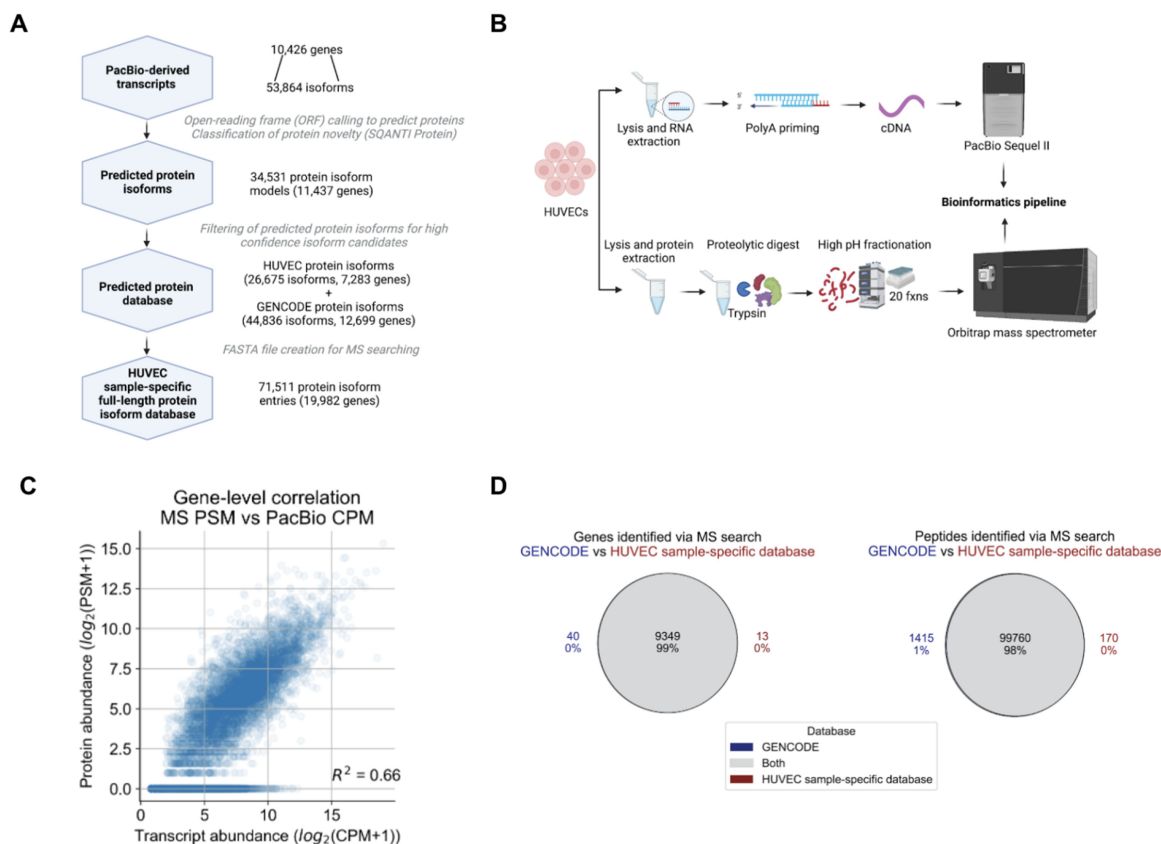


Figure 3. Proteomic analysis of HUVECs using a customized long-read-derived protein isoform database. (A) Steps involved in the generation of a HUVEC sample-specific database. (B) Parallel long-read RNA-seq and MS proteomic data collection from HUVECs. (C) Correlation between estimated RNA and protein expression levels. PSM, peptide spectral match; CPM, full-length read counts per million. (D) Comparison of proteomic search results between the reference and HUVEC sample-specific database.

Table 2. Composition of the HUVEC sample-specific database.

PacBio-derived HUVEC sample-specific database		
Source	Genes	Protein entries
GENCODE	12,699	44,836
PacBio-derived (HUVECs)	7,283	26,675
Contaminants*	-	264
Total	19,982	71,511

*Derived from the MetaMorpheus software version 0.0.316

Information Figure S3). These fractions were then analysed via liquid chromatography LC-MS/MS (Orbitrap Eclipse) in data-dependent acquisition (DDA) mode, generating 3,772,771 MS2 fragmentation spectra. Acquired spectra were searched using the MetaMorpheus [31] software to obtain peptide and protein identifications passing a 1% false-discovery-rate (FDR). Parameters for the MS search can be found in Supplementary Information Table S5.

The HUVEC-specific protein database returns near-complete coverage of detectable peptides from a reference search

To use PacBio-derived transcripts as the basis for deriving a protein database for MS searching, a key assumption is that the detection of a transcript from PacBio data reflects the discovery of a protein product for that isoform, showing that there is

a moderate correlation between transcript and protein abundance. In the past, moderate RNA-protein correlations have been observed using short-read RNA-seq or microarray datasets to quantify transcript abundance [51,52]. Here, we examined the correlation of the transcript abundance that is computed from the long-read RNA-seq data (sum total transcript abundance for a gene, in units of CPM) to the estimated protein abundance (sum total peptide counts passing a 1% FDR, in units of number of peptide spectral matches or PSMs). We observed a moderate correlation with a coefficient of determination (R-square) of 0.66 (Fig. 3C), providing support that the PacBio-based transcript abundances should serve as a reasonable proxy for protein presence, although that may not always be the case for a particular gene.

To assess the general protein sequence content of the HUVEC sample-specific database (not resolved to individual isoforms), we assessed recovery of annotated peptides and genes. The MS data was searched against the GENCODE and UniProt databases to define the set of annotated peptides and genes detectable in the HUVEC sample, and then the same data was searched against the HUVEC sample-specific database. We found that the HUVEC sample-specific database search returned 98% of the peptide and 99% of the gene identifications that were identified when using the GENCODE database for searching (Fig. 3D). The extent of overlap between peptides and genes was similar for the UniProt search results (Supplementary

Information Figure S5). Overall, these results indicate that the HUVEC sample-specific database, which was derived *de novo* from long reads, is able to capture a majority of the detectable gene and peptide populations likely expressed in HUVECs. Confirmation of the large overlap of peptide populations identified by the sample-specific database is ultimately useful since it is the underlying populations of peptides identified that are the basis for protein isoform characterization.

Characterization of HUVEC protein isoforms based on available peptide evidence

We have shown that nearly all reference annotated peptides that are detectable are represented in the HUVEC sample-specific database. With the goal of characterizing isoform expression in endothelial cells, we next evaluated the evidence for the presence of isoforms, in terms of the patterns of their underlying peptide identifications. Due to the complexities and potential ambiguities of protein inference [27], we elected to examine the peptide evidence directly.

We defined three scenarios of isoform detection precision, based on how the set of identified peptides map to isoforms of a gene. The first scenario is when all isoforms of a gene contain only shared peptides, in which the presence of any isoform cannot be definitively confirmed (Fig. 4A, 'Protein isoforms correspond to shared peptides'). Among the 10,444 genes with any peptide evidence, we found that 5,993 genes (57%) were cases in which no isoform could be specifically confirmed as expressed because all mapped peptides were shared among two or more isoforms. Of these genes evidenced only by shared peptides, 3,436 are genes containing PacBio-derived protein isoforms in the hybrid database.

In all other scenarios, there is evidence for the existence of a specific protein isoform because one or more isoforms contain a uniquely mapping peptide. Indeed, the second scenario is when an isoform-specific peptide is identified (Fig. 4A, 'One protein isoform confirmed with a unique peptide'). We found 4,451 (42%) genes for which we have unambiguously identified at least one isoform for a gene. For 1,748 (17%) of genes, only a single isoform was listed in the database, thus, all peptides would be expected to be uniquely mapped. For the remaining 2,703 (26%) of genes with

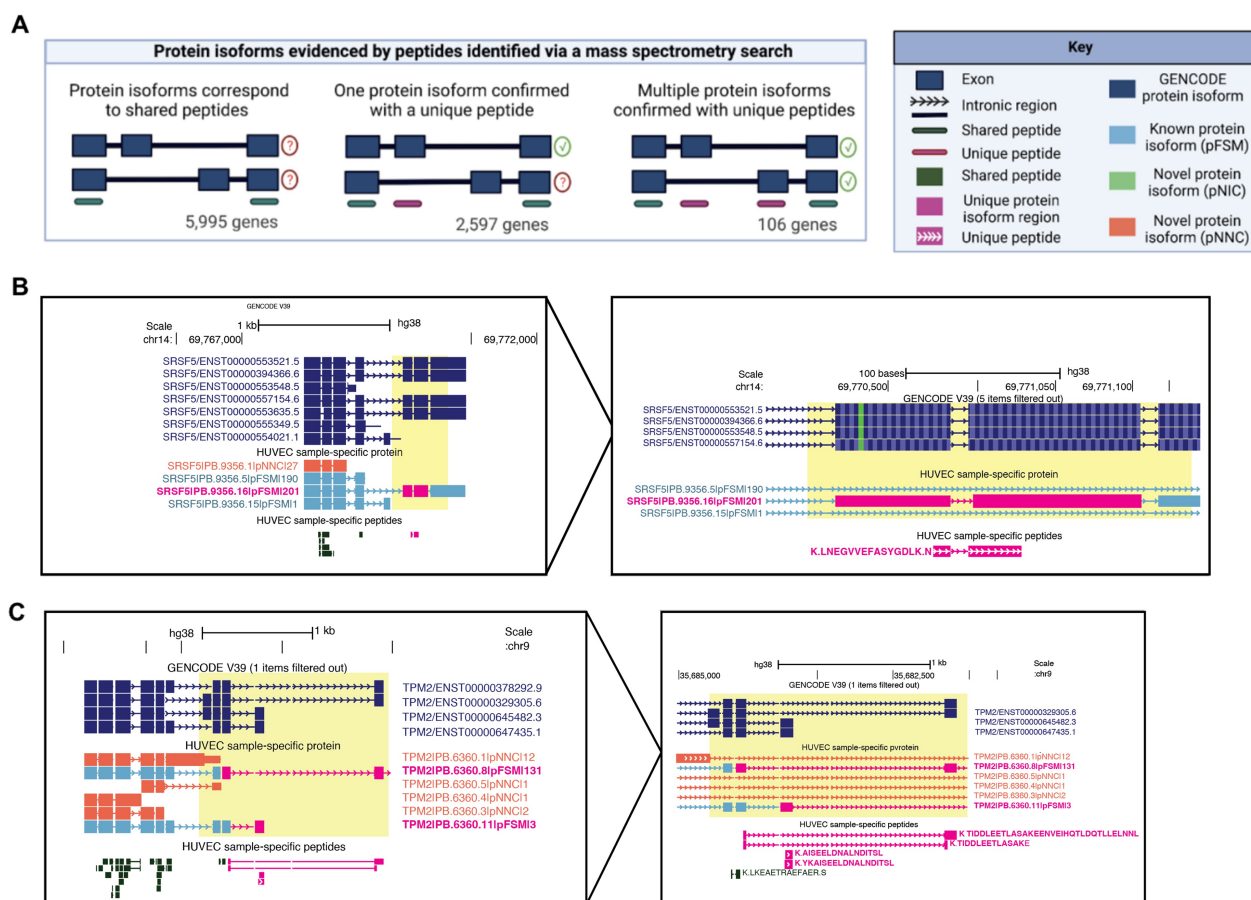


Figure 4. Protein isoforms analyzed based on peptides identified via mass-spectrometry (MS). (A) Scenarios of differing protein isoform detection precision when evidenced by peptides identified from MS. Only genes with multiple protein isoforms in the database are included, and 1,904 genes that express only one isoform were excluded. (B) A protein isoform confirmed with a uniquely mapping peptide LNE, for *SRSF5*, a splice factor that regulates transcripts of VEGF-A. (C) Two protein isoforms of *TPM2* are confirmed with uniquely mapping peptides TID, AIS, and YKA. In B and C PacBio-derived protein isoform label follows this format: <Gene>|<PB accession>|<SQANTI Protein class>|<CPM>.

multiple isoforms annotated, 2,597 (25%) of genes have a single isoform with unique peptide evidence. For example, we found a single isoform supported by a uniquely mapping peptide (Sequence: LNEGVEFASYGDLK) for Serine and Arginine Rich Splicing Factor 5 (*SRSF5*), which is involved in the splicing of VEGF-A pre-mRNA (Fig. 4B). Notably, this peptide is shared among two isoforms in the GENCODE database, meaning that the reference database search results cannot pinpoint the source isoform for this peptide.

Of particular interest is a third scenario in which we found evidence for co-expression of two or more isoforms, each supported by a uniquely mapping peptide. In such cases, a natural question is the nature of the functional relationship between the two isoforms and their biological role in endothelial cells. We found 106 (1%) genes with evidence of two or more co-expressing isoforms (Fig. 4A, 'Multiple protein isoforms confirmed with unique peptides'). For example, we found two isoforms for Tropomyosin 2 (*TPM2*), each supported by a unique peptide (Fig. 4C). Notably, there were nine genes in which three or more isoforms each had unique peptide evidence. Interestingly, there was an unusually large number of seven protein isoforms detected from the gene Plectin (*PLEC*), which exhibits a series of alternative N-termini due to differential 5' transcription (Supporting Information Figure S5). A list of all protein isoforms supported by peptide evidence can be found in Supplementary Information Table S6.

Collectively, these results highlight that while some isoforms may be readily identified from peptide evidence alone, overall, the standard bottom-up MS approach alone does not reach the coverage needed to directly characterize all isoforms predicted from the transcriptome, as observed previously [53,54]. Obtaining peptides suitable to resolve protein isoform identification is limited by the peptides detected during bottom-up MS. Part of the challenge is that when comparing isoforms of the same gene, only small stretches of amino acids are unique to an isoform, while the vast majority of the amino acid sequence is shared [22,53]. Therefore, sampling peptides from the small space of unique amino acids that can directly confirm the presence of a protein isoform is limited by the space of "informative" (i.e. unique-to-an-isoform) peptides[55].

Increased support for protein isoform presence in HUVECs through incorporation of underlying transcript evidence from long-read RNA-seq

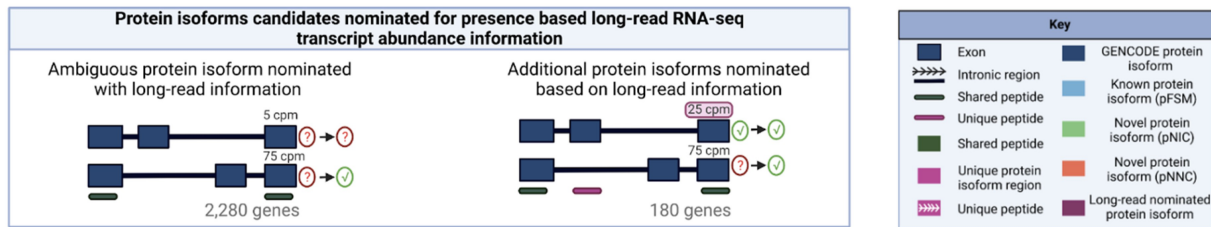
Despite the use of a sample-specific protein isoform database for MS analysis, a large population of predicted isoforms are only supported by shared peptides (Fig. 4A). Because shared peptides map ambiguously to multiple isoforms, they cannot directly confirm expression of any particular isoform in the sample. However, the evidence for a particular protein isoform could be strengthened by considering the underlying transcript abundance levels provided by the sample-matched long-read RNA-seq data, a concept we previously introduced, and has been described for short-read RNA-seq data [22,56,57]. We reasoned that transcript abundance could be used as an additional source of evidence in the isoform discovery process, given there is a moderate correlation between RNA and protein abundance (Fig. 3C).

To explore how long-read RNA-seq data can nominate particular protein isoforms, we first focused on scenarios for which all predicted isoforms for a gene are supported only by shared peptide support. Among such ambiguous protein isoform sets, we reasoned there is higher likelihood for expression of protein isoforms for which the associated transcript abundance is moderately high (e.g. 25 CPM or higher, Fig. 5A). As described in the previous section, 5,993 genes had only shared peptide evidence. Among those genes, 3,436 (57%) contained PacBio-derived isoforms, which have associated transcript abundance information.

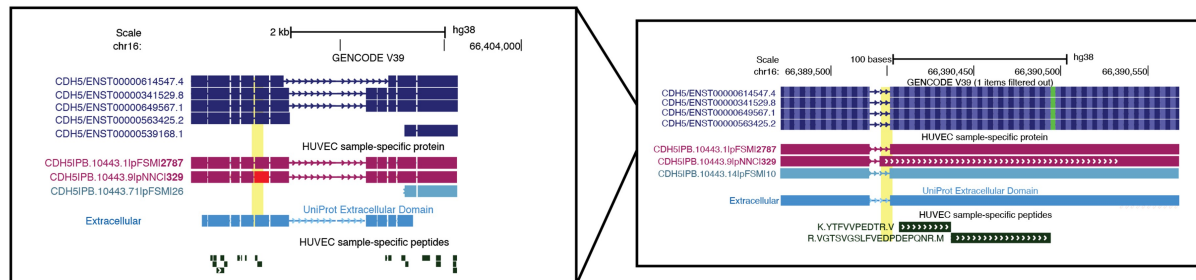
We found that 2,280 (38%) out of the 3,436 genes contain at least one isoform with a moderately high transcript abundance of 25 CPM or higher (Fig. 5A, 'Ambiguous protein isoform nominated with long-read information'). Interestingly, we found 247 (4%) genes in which there is potential co-expression of at least two protein isoforms in HUVECs. For example, we found that *CDH5*, otherwise known as VE-Cadherin [58], potentially expresses multiple protein isoforms. One isoform is highly expressed (PB.10443.1; 2,787 CPM) and matches a protein isoform in GENCODE and UniProt (GENCODE isoform CDH5-201, UniProt accession P33151). However, another isoform is robustly expressed (PB.10443.9, 326 CPM) and, interestingly, is a novel isoform because of alternative usage of a novel splice acceptor (Fig. 5B). This splicing event leads to an isoform of *CDH5* that gains nine amino acids in the extracellular domain, the region of the protein responsible for mediating interactions with other cadherins to regulate endothelial adhesion properties. This example highlights that while *CDH5* isoforms were only supported by shared peptides, the incorporation of the transcript abundance information as provided by the matched long-read data provides higher weights on the existence of at least two isoforms.

To further explore how long-read RNA-seq data can provide additional evidence for expression of protein isoforms, we focused on scenarios in which there is clear evidence for one isoform based on unique peptide evidence, but another isoform of the same gene is supported by only shared peptides (Fig. 5A, 'Additional protein isoforms nominated based on long-read information'). We found 180 genes (3%) for which the existence of the alternative protein isoform is supported by long-read evidence (i.e. 25 CPM or higher transcript abundance). Interestingly, we found several protein isoforms of a key endothelial cell surface marker, the platelet endothelial cell adhesion molecule, *PECAM1* (also known as *CD31*) [59]. We found a unique peptide identified for *PECAM1* (PB.11293.25, Sequence: SDSGTYICTAEMLSQPR), but the remainder of peptides identified for *PECAM1* are shared across multiple *PECAM1* isoforms, leaving open uncertainty about the expression of other *PECAM1* isoforms beyond PB.11293.25. From the transcript abundance information, we nominated three additional isoforms accompanied by strong long-read support for *PECAM1* (PB 11293.22, 75 CPM; PB 11293.1, 79 CPM; PB 11293.7, 543 CPM; Fig. 5C). *PECAM1* produces a transmembrane protein with an extracellular domain, transmembrane-spanning domain, and a C-terminal cytoplasmic domain that likely interacts with intracellular signalling proteins in endothelial cells [59–61]. Strikingly, the differential exon usage observed for these three isoforms are located exclusively in the C-terminal domain, suggesting

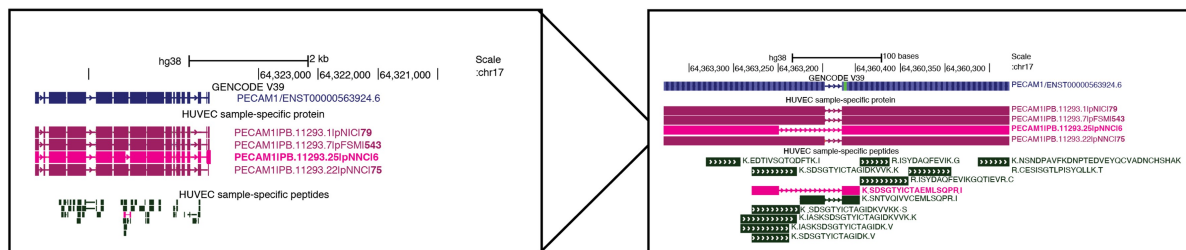
A



B



C



*Not all protein isoforms shown

Figure 5. Nomination of protein isoforms when incorporating long-read data. (A) Scenarios of protein isoform candidates nominated for expression when transcript abundance from the long-read RNA-seq information is incorporated. (B) *CDHS* gene, involved in endothelial pathways demonstrating a scenario of ambiguous protein isoforms identified only by shared peptides, but incorporation of long-read RNA-seq data suggests the expression of three moderately expressed protein isoforms (PB.10443.1, PB.10443.9 and PB.10443.71). (C) *PECAM1* gene, involved in endothelial pathways demonstrating an example where one protein isoform is identified via a unique peptide (PB.11293.25), SDS, while the remaining protein isoforms are supported by shared peptides. Abundance information from long-read RNA-seq suggest expression of (PB.11293.1 and PB.11293.7). In B and C, PacBio-derived protein isoform label follows this format: <Gene>|<PB accession>|<SQANTI Protein class>|<CPM>. For B and C, low abundance protein isoforms (<25 CPM) are not shown.

potential changes to interactions with intracellular signalling molecules. Further details on candidates identified via long-read abundance information can be found in Supplementary Information Table S6.

Collectively, these case studies highlight how incorporation of transcript abundance information could nominate protein isoforms which were unable to be directly confirmed as expressed based solely on MS peptide evidence. Note that this approach does not provide any information on the absence of protein isoforms with lower transcript abundance, but, rather, is supplying additional lines of evidence to nominate protein isoforms that may have higher likelihood of expression and represent candidates for functional study. Such isoforms are attractive candidates for further MS validation and subsequent functional analysis.

Novel protein isoform discovery enabled through the HUVEC sample-specific database

We have shown that utilization of a HUVEC sample-specific protein database, with the accompanying transcript abundance

values, can lead to inference of novel protein isoform presence. A more direct way to confirm the presence of a novel protein isoform is by detecting a uniquely mapping novel peptide. However, the knowledge of the full-length protein isoforms expressed within a sample is not always possible when using short-read RNA-seq, which can return information on individual splice junctions but may not accurately define full-length transcripts [17]. Long-read RNA-seq provides the full-length transcript and, by extension, the full-length protein isoform prediction; therefore, a novel peptide that directly maps to the full-length protein isoform lends support for its existence.

Using the sample-specific database, we discovered novel peptides for HUVECs, indicating that the reference proteome does not comprehensively capture all protein isoform diversity in a sample. We found 108 novel peptide sequences passing a global 1% FDR, for which they are not represented within the GENECode or UniProt databases (Fig. 6A, Supplementary Information Table S7) [35,36]. Increased false positive rates for novel peptides have been observed previously [62]; therefore, we employed strict validation criteria for the novel peptides. Of the 108 novel peptides identified, 39 peptides had a Q-value score

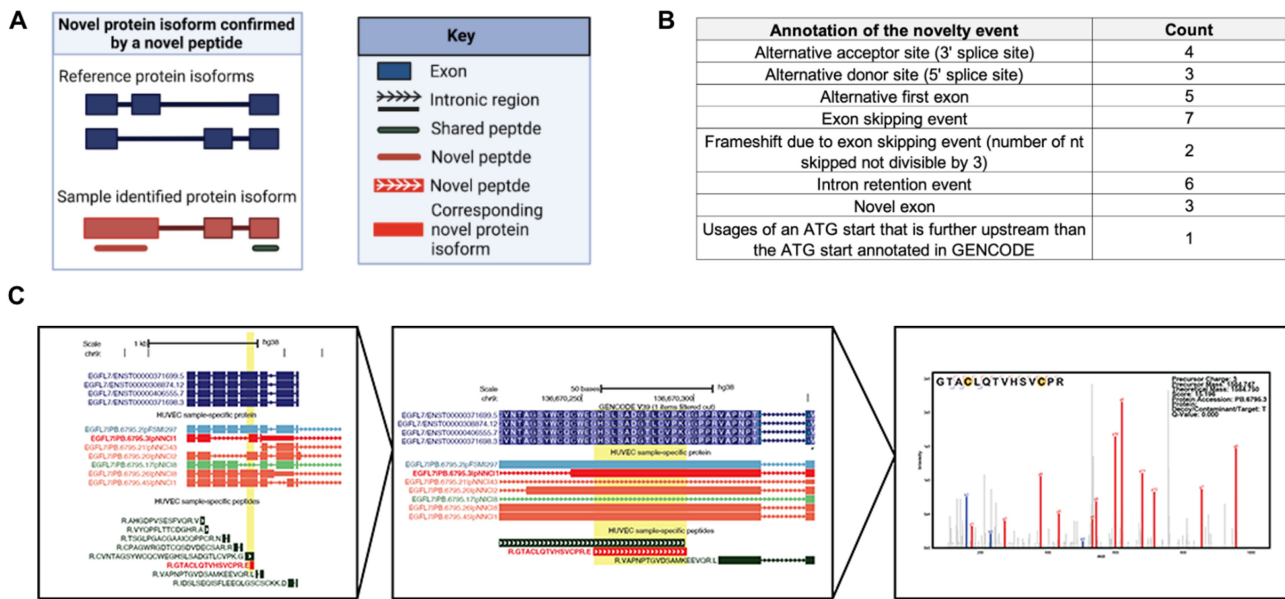


Figure 6. Novel protein isoforms discovered via unique peptides. (A) Novel protein isoform confirmed by identified novel peptides. (B) Table of the frequency of events supported confirmation of a novel peptide. (C) Novel peptide found for a protein isoform of endothelial gene *EGFL7*. Novel peptide and corresponding protein isoform shown in red, which supports a frameshift event for the protein isoform PB.6795.3.

below 0.001, corresponding to a 0.1% global FDR. Upon manual annotation of these 39 peptides, we noted 30 peptides with especially strong spectral support, such as full ladders of b- and y- ion fragmentation peaks in the MS2 raw spectra. These novel peptides supported expression of novel alternative transcription or splicing events, such as retained intronic regions or novel exons (Fig. 6B, Supplementary Information Table S7).

Of the identified novel isoforms, we closely examined splicing events for genes previously implicated in endothelial pathways. Such novel isoforms could represent attractive candidate isoforms for further functional characterization. We found a novel peptide (Sequence: GTACLQTVHVSVCPR) confirming the expression of a splice-induced frame-shifted region of *EGFL7* (Protein entry: PB.6795.3), a gene reported through the literature to be involved in vasculogenic pathways as well as hemogenic specification (Fig. 6C) [63,64]. We also discovered two novel peptides for *PECAM1*. This is an important finding since *PECAM1* is a marker for endothelial cells and plays a role in the regulation of junctional integrity of endothelial cells and vascular barrier [59]. Specifically, we discovered a novel peptide (Sequence: ELELLTSKDPSPASQSAGITDLGKK, maps to protein entry PB.11293.45) corresponding to a novel exon, as well as a second novel peptide (Sequence: SDSGTICTAEMLSQPR, mapping to protein entry PB.11293.25) that confirms the usage of a novel alternative donor site (Supporting information Table S6).

Conclusions

Endothelial cells that line all blood vessels are critical for the cardiovascular system and their behaviours can be modulated by protein isoforms, though the extent of this mechanism is not known. To characterize isoform expression in endothelial cells, we performed long-read RNA sequencing (PacBio) of HUVECs to characterize transcript isoforms, and predicted proteins via their translation *in silico* to protein isoform

sequences. To assess evidence for protein isoform expression, we performed MS analysis on the same HUVEC sample and used the HUVEC sample-specific database for MS searching. This general approach has been described and termed ‘long-read proteogenomics’ to enhance protein isoform characterization [22].

Our long-read proteogenomics workflow applied to HUVECs, led to the identification of 53,863 distinct transcript isoforms, of which 22,195 were novel. We also found 8,522 genes co-expressing multiple isoforms. Surprisingly, a quarter of the time, the most abundant isoform in HUVECs did not match the predicted ‘reference isoform’ (GENCODE APPRIS principal isoform). This includes genes annotated in endothelial pathways including *CD34* and *NRP1*. From the transcript sequences, we derived a hybrid protein isoform database that contains the highest confidence protein isoform predictions from PacBio-derived transcript isoform sequences, which was completed with GENCODE reference and contaminant sequences. The long-read-derived database captures almost all peptides and proteins detected from searches against the GENCODE protein database.

We identified 10,444 genes with peptide evidence. Based on the peptides identified through MS searching, we found support for expression of 4,451 genes based on uniquely mapping peptides. For the remaining 5,993 genes only evidenced by shared peptides, we incorporated the underlying transcript abundance information as an additional layer of evidence, nominating an additional 2,280 genes as potentially expressed. This group includes a novel isoform for endothelial gene *CDH5* (VE-Cadherin). This case exemplifies how a combination of the full-length transcript and proteomics data can lead to the discovery of novel protein isoforms that cannot be identified by MS data alone. We showed that the HUVEC sample-specific database enabled discovery of 108 novel protein isoforms based on novel peptide identifications.

Among the novel protein isoforms identified is the endothelial gene *PECAM1*.

Our proteogenomic method shows promise for isoform discovery in endothelial cells, but opportunities exist for further improvements. First, limitations in the MS coverage mean that proteins with low abundance or poorly ionizable peptides remain undetected. Future work could involve targeted proteomics, such as parallel reaction monitoring or advanced targeted acquisition strategies, for sensitive detection of alternative protein isoforms [65–67]. Second, the isoforms discovered in this study represent the results of a single cell line in a static culture condition. For the purposes of identifying isoforms that are dynamically regulated, multiple conditions should be examined. Third, the sample-specific database relies on the assumption that sequenced transcripts reflect protein sequences. Thus, we assume that transcripts are both fully sampled as well as moderately correlated to protein expression, which may not be the case for all genes. And finally, our pipeline so far is focused on proteins arising from genes already annotated as protein-coding. An interesting future direction would be to include long non-coding RNAs or other ostensibly non-coding transcripts, which may reveal coding potential through the proteogenomics approach [68].

Overall, we have shown the application of a long-read proteogenomics platform towards characterization of known and novel isoforms in primary endothelial cells. This approach can uncover isoform populations that could modulate endothelial cell phenotype and function. The systematic discovery of isoforms produces information to guide selection of candidate isoforms for functional studies. This approach can be extended to various endothelial cell contexts including both healthy and diseased states to chart isoforms changing across development or during onset of cardiovascular disease.

Acknowledgments

This work was financially supported by the Robert M. Berne Cardiovascular Research Center Training Program (T32HL007284) to MMM. We gratefully acknowledge the additional support by the NIH HL146056 and DK118728 to KKH. The long-read sequencing was performed at the Maryland Genomics at the University of Maryland Institute for Genome Science. Figs. 1–6 were created with Biorender.com. The graphical abstract was created with the help from Katharine Tuttle. We thank Dr. Nicholas Chavkin and Jordon Aragon for helpful comments and feedback related to this project.

Author Contributions

GMS and MMM conceived of the project. GMS designed the study and supervised the project, along with KKH. MMM was involved in data collection, data analysis, interpretation, and conclusions, with discussions with GMS. LS performed the long-read RNA-seq experimental analysis. BTJ performed computational analysis, including long-read RNA-seq and proteomics analysis. MMM and EDJ performed the novel peptide analysis. JS and MM performed computational analysis of novel isoforms. BTJ, MMM and GMS contributed to analysis reproducibility, data curation and design of the workflow for the data described in this paper. MMM and GMS wrote the manuscript with input from all authors. All authors read and approved the final manuscript.

Disclosure statement

No potential conflict of interest was reported by the authors.

Funding

This work was supported by the National Heart, Lung, and Blood Institute (NHLBI) [HL146056]; National Institute of Diabetes and Digestive and Kidney Diseases [DK118728]; National Heart, Lung, and Blood Institute (NHLBI) [T32HL007284].

Abbreviations

Alternative splicing (AS)
 Annotation of principal and alternative splice isoforms (APPRIS)
 Circular consensus reads (CCS)
 Counts per million (CPM)
 Enriched GO terms (GORilla)
 False-discovery-rate (FDR)
 Fast-all (FASTA)
 Filter-aided sample preparation (FASP)
 Full splice match (FSM)
 Human umbilical vein endothelial cells (HUVEC)
 Iso-seq (PacBio processing)
 Mass-spectrometry (MS)
 Messenger RNA (mRNA)
 Novel in catalogue (NIC)
 Novel not in catalogue (NNC)
 Open-reading frame (ORF)
 RNA-sequencing (RNA-seq)
 Protein full splice match (pFSM)
 Protein novel in catalogue (pNIC)
 Protein novel not in catalogue (pNNC)
 Single-Molecule Real-Time sequencing (SMRT-seq)
 Serine and Arginine Rich Splicing Factor 5 (SRSF5)
 Structural and Quality Annotation of Novel Transcript Isoforms (SQANTI)
 Vascular endothelial growth factor A (VEGF-A)

ORCID

Gloria M. Sheynkman  <http://orcid.org/0000-0002-4223-9947>

References

- [1] Cleaver O, Melton DA. Endothelial signaling during development. *Nat Med* [Internet] 2003; 9:661–668. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/12778164>
- [2] Rajendran P, Rengarajan T, Thangavel J, et al. The vascular endothelium and human diseases. *Int J Biol Sci* [Internet] 2013 cited 2021 Mar 23; 9:1057–1069. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3831119>
- [3] Richardson MR, Lai X, Witzmann FA, et al. Venous and arterial endothelial proteomics: mining for markers and mechanisms of endothelial diversity. *Expert Rev Proteomics* [Internet] 2010; 7:823–831. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/21142885>
- [4] Nordon I, Brar R, Hinchliffe R, et al. The role of proteomic research in vascular disease. *J Vasc Surg* [Internet] 2009; 49:1602–1612. Available from: <http://www.sciencedirect.com/science/article/pii/S0741521409006053>
- [5] Farrokh S, Brillen AL, Haendeler J, et al. Critical regulators of endothelial cell functions: for a change being alternative. *Antioxid Redox Signal* [Internet] 2015; 22:1212–1229. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/25203279>
- [6] Bowler E, Oltean S. Alternative splicing in angiogenesis. *Int J Mol Sci* [Internet] 2019; 20:2067. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/31027366>

- [7] Mthembu NN, Mbita Z, Hull R, et al. Abnormalities in alternative splicing of angiogenesis-related genes and their role in HIV-related cancers. *HIV AIDS* [Internet] 2017; 9:77–93.
- [8] Murphy PA, Butty VL, Boutz PL, et al. Alternative RNA splicing in the endothelium mediated in part by Rbfox2 regulates the arterial response to low flow. *Elife* [Internet] 2018; 7. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/29293084>
- [9] Di Matteo A, Belloni E, Pradella D, et al. Alternative splicing in endothelial cells: novel therapeutic opportunities in cancer angiogenesis. *J Exp Clin Cancer Res* [Internet] 2020; 39:275. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/33287867>
- [10] Hang X, Li P, Li Z, et al. Transcription and splicing regulation in human umbilical vein endothelial cells under hypoxic stress conditions by exon array. *BMC Genomics*. 2009;10:126.
- [11] Giampietro C, Deflorian G, Gallo S, et al. The alternative splicing factor Nova2 regulates vascular development and lumen formation. *Nat Commun* [Internet] 2015; 6:8479. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/26446569>
- [12] Khan S, Taverna F, Rohlenova K, et al. EndoDB: a database of endothelial cell transcriptomics data. *Nucleic Acids Res*. 2019;47:D736–44.
- [13] Mudge JM, Harrow J. The state of play in higher eukaryote gene annotation. *Nat Rev Genet* [Internet]. 2016;17:758–772.
- [14] Caniuguir A, Krause BJ, Hernandez C, et al. Markers of early endothelial dysfunction in intrauterine growth restriction-derived human umbilical vein endothelial cells revealed by 2D-DIGE and mass spectrometry analyses. *Placenta* [Internet]. 2016; 41:14–26. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/27208404>
- [15] Banarjee R, Sharma A, Bai S, et al. Proteomic study of endothelial dysfunction induced by AGEs and its possible role in diabetic cardiovascular complications. *J Proteomics* [Internet] 2018; 187:69–79. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/29935336>
- [16] Madugundu AK, Na CH, Nirujogi RS, et al. Integrated transcriptomic and proteomic analysis of primary human umbilical vein endothelial cells. *Proteomics* [Internet] 2019; 19:e1800315. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/30983154>
- [17] Steijger T, Abril JF, Engstrom PG, et al. Assessment of transcript reconstruction methods for RNA-seq. *Nat Methods* [Internet] 2013; 10:1177–1184. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/24185837>
- [18] Nesvizhskii AI. Proteogenomics: concepts, applications and computational strategies. *Nat Methods* [Internet]. 2014;11:1114–1125.
- [19] van Dijk EL, Jaszczyszyn Y, Naquin D, et al. The third revolution in sequencing technology. *Trends Genet* [Internet]. 2018;34:666–681.
- [20] Eid J, Fehr A, Gray J, et al. Real-time DNA sequencing from single polymerase molecules. *Science*. 2009;323:133–138.
- [21] Rhoads A, Au KF. PacBio sequencing and its applications. *Genomics Proteomics Bioinformatics* [Internet] 2015; 13:278–289. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/26542840>
- [22] Miller RM, Jordan BT, Mehlferber MM, et al. Enhanced protein isoform characterization through long-read proteogenomics. *Genome Biol* [Internet] 2022 cited 2022 Mar 4; 23:1–28. Available from: <https://genomebiology.biomedcentral.com/articles/10.1186/s13059-022-02624-y>
- [23] Sharon D, Tilgner H, Grubert F, et al. A single-molecule long-read survey of the human transcriptome. *Nat Biotechnol* [Internet]. 2013;31:1009–1014
- [24] Deslattes Mays A, Schmidt M, Graham G, et al. Single-Molecule Real-Time (SMRT) full-length RNA-sequencing reveals novel and distinct mRNA isoforms in human bone marrow cell subpopulations. *Genes (Basel)* [Internet] 2019; 10:17. Available from: https://res.mdpi.com/d_attachment/genes/genes-10-00253/article_deploy/genes-10-00253-v2.pdf
- [25] Verbruggen S, Gessulat S, Gabriels R, et al. Spectral prediction features as a solution for the search space size problem in proteogenomics. *Mol Cell Proteomics* [Internet]. 2021;2021(04/07):100076. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/33823297>
- [26] Anvar SY, Allard G, Tseng E, et al. Full-length mRNA sequencing uncovers a widespread coupling between transcription initiation and mRNA processing [Internet]. *Genome Biol*. 2018;19. DOI:10.1186/s13059-018-1418-0.
- [27] Nesvizhskii AI, Aebersold R. Interpretation of shotgun proteomic data. *Mol Cell Proteomics* [Internet]. 2005 [cited 2021 Jul 28];4:1419–1440. Available from: [https://www.mcponline.org/article/S1535-9476\(20\)30061-X/abstract](https://www.mcponline.org/article/S1535-9476(20)30061-X/abstract)
- [28] Wiśniewski JR. Filter-aided sample preparation for proteome analysis [Internet]. *Methods Mol Biol*. 2018;3–10. DOI:10.1007/978-1-4939-8695-8_1.
- [29] Tardaguila M, de la Fuente L, Marti C, et al. SQANTI: extensive characterization of long-read transcript sequences for quality control in full-length transcriptome identification and quantification. *Genome Res* [Internet] 2018; 2018 February 15. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/29440222>
- [30] Wang L, Park HJ, Dasari S, et al. CPAT: coding-potential assessment tool using an alignment-free logistic regression model. *Nucleic Acids Res* [Internet]. 2013;41:e74. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/23335781>
- [31] Solntsev SK, Shortreed MR, Frey BL, et al. Enhanced global post-translational modification discovery with metamorphus. *J Proteome Res* [Internet]. 2018;17:1844–1851. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/29578715>
- [32] Elias JE, Gygi SP. Target-decoy search strategy for increased confidence in large-scale protein identifications by mass spectrometry. *Nat Methods* [Internet]. 2007;4:207–214. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/17327847>
- [33] Smith LM, Agar JN, Chamot-Rooke J, et al. Consortium for top-down proteomics. the human proteoform project: defining the human proteome. *Sci Adv* [Internet];. Available from. 2021;7:eabk0734.
- [34] Gordon SP, Tseng E, Salamov A, et al. Widespread polycistronic transcripts in fungi revealed by single-molecule mRNA sequencing. *PLoS One*. 2015;10. Available from <https://journals.plos.org/plosone/article/file?id=10.1371/journal.pone.0132628&type=printable>
- [35] Frankish A, Diekhans M, Ferreira AM, et al. GENCODE reference annotation for the human and mouse genomes. *Nucleic Acids Res* [Internet]. 2019;47:D766–73. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6323946/pdf/gky955.pdf>
- [36] UniProt C. UniProt: a worldwide hub of protein knowledge. *Nucleic Acids Res* [Internet]. 2019;47: D506–15. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/30395287>
- [37] McGarvey PB, Nightingale A, Luo J, et al. UniProt genomic mapping for deciphering functional effects of missense variants. *Hum Mutat* [Internet] 2019; 40:694–705. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/30840782>
- [38] Huang KK, Huang J, Jkl W, et al. Long-read transcriptome sequencing reveals abundant promoter diversity in distinct molecular subtypes of gastric cancer. *Genome Biol* [Internet] 2021; 22:44. Available from:
- [39] Leung SK, Jeffries AR, Castanho I, et al. Full-length transcript sequencing of human and mouse cerebral cortex identifies widespread isoform diversity and alternative splicing. *Cell Rep* [Internet]. 2021;37:110022.
- [40] Marcelo KL, Goldie LC, Hirschi KK. Regulation of endothelial cell differentiation and specification. *Circ Res* [Internet] 2013; 112:1272–1287. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/23620236>
- [41] Aragon JW, Hirschi KK. Endothelial Cell Differentiation and Hemogenic Specification. *Cold Spring Harb Perspect Med* [Internet] 2022; Available from:;12(7):a041164.
- [42] Eden E, Navon R, Steinfeld I, et al. GOrilla: a tool for discovery and visualization of enriched GO terms in ranked gene lists. *BMC Bioinformatics* [Internet] 2009; 10:48.
- [43] Lanahan A, Zhang X, Fantin A, et al. The neuropilin 1 cytoplasmic domain is required for VEGF-A-dependent arteriogenesis. *Dev Cell* [Internet]. 2013;25:156–168.
- [44] Kofler NM, Simons M. Angiogenesis versus arteriogenesis: neuropilin 1 modulation of VEGF signaling. *F1000Prime Rep* [Internet] 2015; 7:26.
- [45] Chang K-T, Wang L-H, Lin Y-M, et al. CELF1 promotes vascular endothelial growth factor degradation resulting in impaired microvasculature in heart failure. *FASEB J* [Internet]. 2021;35:e21512.

- [46] Edwards JM, Long J, de Moor Ch, et al. Structural insights into the targeting of mRNA GU-rich elements by the three RRM domains of CELF1. *Nucleic Acids Res* [Internet]. 2013;41:7153–7166.
- [47] Blech-Hermoni Y, Stillwagon SJ, Ladd AN. Diversity and conservation of CELF1 and CELF2 RNA and protein expression patterns during embryonic development. *Dev Dyn* [Internet]. 2013;242:767–777. Available from: <https://anatomypubs.onlinelibrary.wiley.com/doi/abs/10.1002/dvdy.23959>
- [48] Rodriguez JM, Rodriguez-Rivas J, Di Domenico T, et al. APPRIS 2017: principal isoforms for multiple gene sets. *Nucleic Acids Res* [Internet]. 2018;46:D213–7.
- [49] Wang X, Slebos RJ, Wang D, et al. Protein identification using customized protein sequence databases derived from RNA-Seq data. *J Proteome Res* [Internet]. 2012 cited 2012 Jan 5;11:1009–1017.
- [50] Van Nostrand EL, Freese P, Pratt GA, et al. A large-scale binding and functional map of human RNA-binding proteins. *Nat*. 2020;583:711–719.
- [51] Wang D, Eraslan B, Wieland T, et al. A deep proteome and transcriptome abundance atlas of 29 healthy human tissues. *Mol Syst Biol* [Internet]. 2019; 15:e8503. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/30777892>
- [52] Vogel C, Marcotte EM. Insights into the regulation of protein abundance from proteomic and transcriptomic analyses. *Nat Rev Genet* [Internet]. 2012;13:227–232.
- [53] Blakeley P, Siepen JA, Lawless C, et al. Investigating protein isoforms via proteomics: a feasibility study. *Proteomics* [Internet]. 2010; 10:1127–1140.
- [54] Lau E, Han Y, Williams DR, et al. Splice-junction-based mapping of alternative isoforms in the human proteome. *Cell Rep* [Internet]. 2019 [cited 2020 Feb 20];29:3751–65.e5. Available from: <https://linkinghub.elsevier.com/retrieve/pii/S2211124719314974>
- [55] Wang X, Codreanu SG, Wen B, et al. Detection of proteome diversity resulted from alternative splicing is limited by trypsin cleavage specificity. *Mol Cell Proteomics* [Internet]. 2018;17:422–430.
- [56] Carlyle BC, Kitchen RR, Zhang J, et al. Isoform-level interpretation of high-throughput proteomics data enabled by deep integration with RNA-seq. *J Proteome Res* [Internet]. 2018;17:3431–3444. Available from: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6392456/pdf/nihms-1012113.pdf>
- [57] Salovska B, Zhu H, Gandhi T, et al. Isoform-resolved correlation analysis between mRNA abundance regulation and protein level degradation. *Mol Syst Biol* [Internet]. 2020; 16:e9170. Available from: <https://www.embopress.org/doi/abs/10.15252/msb.20199170>
- [58] Sauter L, Krudewig A, Herwig L, et al. Cdh5/VE-cadherin promotes endothelial cell interface elongation via cortical actin polymerization during angiogenic sprouting. *Cell Rep* [Internet]. 2014;9:504–513.
- [59] Privratsky JR, Newman PJ. PECAM-1: regulator of endothelial junctional integrity. *Cell Tissue Res* [Internet]. 2014;355:607–619.
- [60] Cao G, O'Brien CD, Zhou Z, et al. Involvement of human PECAM-1 in angiogenesis and in vitro endothelial cell migration. *Am J Physiol Cell Physiol* [Internet]. 2002; 282:C1181–90.
- [61] Dusserre N, L'Heureux N, Bell KS, et al. PECAM-1 interacts with nitric oxide synthase in human endothelial cells: implication for flow-induced nitric oxide synthase activation. *Arterioscler Thromb Vasc Biol* [Internet]. 2004;24:1796–1802.
- [62] Castellana N, Bafna V. Proteogenomics to discover the full coding content of genomes: a computational perspective. *J Proteomics* [Internet]. 2010; 73:2124–2135. Available from: <https://www.ncbi.nlm.nih.gov/pubmed/20620248>
- [63] Nichol D, Stuhlmann H. EGFL7: a unique angiogenic signaling factor in vascular development and disease. *Blood*. 2012;119:1345–1352.
- [64] Schmidt MHH, Bicker F, Nikolic I, et al. Epidermal growth factor-like domain 7 (EGFL7) modulates Notch signalling and affects neural stem cell renewal. *Nat Cell Biol* [Internet]. 2009;11:873–880.
- [65] Gallien S, Kim SY, Domon B. Large-scale targeted proteomics using internal standard triggered-parallel reaction monitoring (IS-PRM)*[S]. *Mol Cell Proteomics* [Internet]. 2015;14:1630–1644. Available from: [https://www.mcponline.org/article/S1535-9476\(20\)33173-X/abstract](https://www.mcponline.org/article/S1535-9476(20)33173-X/abstract)
- [66] Erickson BK, Rose CM, Braun CR, et al. A strategy to combine sample multiplexing with targeted proteomics assays for high-throughput protein signature characterization. *Mol Cell* [Internet]. 2017;65:361–370.
- [67] Wichmann C, Meier F, Winter SV, et al. MaxQuant.Live enables global targeting of more than 25,000 peptides. *Mol Cell Proteomics*. 2019;18(5):982–994. DOI:10.1074/mcp.TIR118.001131.
- [68] Mattick JS. 2018. The state of long non-coding RNA biology. *Noncoding RNA* [Internet]; 4. DOI:10.3390/ncrna4030017.