

# UC San Diego

## UC San Diego Previously Published Works

### Title

Kawasaki Disease Patient Stratification and Pathway Analysis Based on Host Transcriptomic and Proteomic Profiles

### Permalink

<https://escholarship.org/uc/item/6m23886f>

### Journal

International Journal of Molecular Sciences, 22(11)

### ISSN

1661-6596

### Authors

Jackson, Heather  
Menikou, Stephanie  
Hamilton, Shea  
et al.

### Publication Date

2021

### DOI

10.3390/ijms22115655

Peer reviewed



Article

# Kawasaki Disease Patient Stratification and Pathway Analysis Based on Host Transcriptomic and Proteomic Profiles

Heather Jackson <sup>1,†</sup>, Stephanie Menikou <sup>1,†</sup>, Shea Hamilton <sup>1</sup>, Andrew McArdle <sup>1</sup>, Chisato Shimizu <sup>2</sup>, Rachel Galassini <sup>1</sup>, Honglei Huang <sup>3</sup>, Jihoon Kim <sup>4</sup>, Adriana Tremoulet <sup>2</sup>, Adam Thorne <sup>5</sup>, Roman Fischer <sup>6</sup>, Marien I. de Jonge <sup>7</sup>, Taco Kuijpers <sup>8</sup>, Victoria Wright <sup>1</sup>, Jane C. Burns <sup>2</sup>, Climent Casals-Pascual <sup>9</sup>, Jethro Herberg <sup>1</sup>, Mike Levin <sup>1,‡</sup>, Myrsini Kaforou <sup>1,\*</sup> and on behalf of the PERFORM Consortium <sup>§</sup>

<sup>1</sup> Faculty of Medicine, Imperial College London, London SW7 2AZ, UK; heather.jackson17@imperial.ac.uk (H.J.); s.menikou@imperial.ac.uk (S.M.); s.hamilton@imperial.ac.uk (S.H.); a.mcardle@imperial.ac.uk (A.M.); r.galassini@imperial.ac.uk (R.G.); v.wright@imperial.ac.uk (V.W.); j.herberg@imperial.ac.uk (J.H.); m.levin@imperial.ac.uk (M.L.)

<sup>2</sup> Department of Pediatrics, University of California San Diego, La Jolla, CA 92093, USA; c1shimizu@health.ucsd.edu (C.S.); atremoulet@health.ucsd.edu (A.T.); jcburns@health.ucsd.edu (J.C.B.)

<sup>3</sup> Target Discovery Institute, University of Oxford, Oxford OX3 7FZ, UK; hlhuang318@hotmail.com

<sup>4</sup> Department of Biomedical Informatics, University of California San Diego, La Jolla, CA 92093, USA; j5kim@health.ucsd.edu

<sup>5</sup> Department of Surgery, Section of Hepatobiliary Surgery and Liver Transplantation, University of Groningen, University Medical Center Groningen, 9713 GZ Groningen, The Netherlands; a.m.thorne@umcg.nl

<sup>6</sup> Discovery Proteomics Facility, Target Discovery Institute, Nuffield Department of Medicine, University of Oxford, Oxford OX3 9DU, UK; roman.fischer@ndm.ox.ac.uk

<sup>7</sup> Section Pediatric Infectious Diseases, Laboratory of Medical Immunology, Radboud Institute for Molecular Life Sciences, Radboud University Medical Center, 6525 GA Nijmegen, The Netherlands; marien.dejonge@radboudumc.nl

<sup>8</sup> Department of Pediatric Immunology, Rheumatology, and Infectious Diseases, Emma Children's Hospital, Amsterdam University Medical Center (AMC), 1105 AZ Amsterdam, The Netherlands; t.w.kuijpers@amsterdamumc.nl

<sup>9</sup> Wellcome Trust Centre for Human Genetics, University of Oxford, Oxford OX3 7BN, UK; ccasals@clinic.cat

\* Correspondence: m.kaforou@imperial.ac.uk

† Shared first author.

‡ Shared final author.

§ Personalised Risk Assessment in Febrile Illness to Optimise Real-Life Management (PERFORM), London SW7 2AZ, UK.



**Citation:** Jackson, H.; Menikou, S.; Hamilton, S.; McArdle, A.; Shimizu, C.; Galassini, R.; Huang, H.; Kim, J.; Tremoulet, A.; Thorne, A.; et al. Kawasaki Disease Patient Stratification and Pathway Analysis Based on Host Transcriptomic and Proteomic Profiles. *Int. J. Mol. Sci.* **2021**, *22*, 5655. <https://doi.org/10.3390/ijms22115655>

Academic Editor: Sarath Janga Chandra

Received: 9 April 2021

Accepted: 4 May 2021

Published: 26 May 2021

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



**Copyright:** © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

**Abstract:** The aetiology of Kawasaki disease (KD), an acute inflammatory disorder of childhood, remains unknown despite various triggers of KD having been proposed. Host 'omic profiles offer insights into the host response to infection and inflammation, with the interrogation of multiple 'omic levels in parallel providing a more comprehensive picture. We used differential abundance analysis, pathway analysis, clustering, and classification techniques to explore whether the host response in KD is more similar to the response to bacterial or viral infections at the transcriptomic and proteomic levels through comparison of 'omic profiles from children with KD to those with bacterial and viral infections. Pathways activated in patients with KD included those involved in anti-viral and anti-bacterial responses. Unsupervised clustering showed that the majority of KD patients clustered with bacterial patients on both 'omic levels, whilst application of diagnostic signatures specific for bacterial and viral infections revealed that many transcriptomic KD samples had low probabilities of having bacterial or viral infections, suggesting that KD may be triggered by a different process not typical of either common bacterial or viral infections. Clustering based on the transcriptomic and proteomic responses during KD revealed three clusters of KD patients on both 'omic levels, suggesting heterogeneity within the inflammatory response during KD. The observed heterogeneity may reflect differences in the host response to a common trigger, or variation dependent on different triggers of the condition.

**Keywords:** infectious diseases; paediatrics; transcriptomics; proteomics; Kawasaki disease; host 'omics; systems biology; pathway analysis; clustering; classification

## 1. Introduction

Kawasaki disease (KD) is an acute inflammatory disorder first described in Japan over 50 years ago [1]. KD occurs most frequently in children under five years of age [2]. Untreated KD leads to the formation of coronary artery aneurysms (CAAs) in 10–30% of children [3–5], making it the most common cause of acquired heart disease in children in Europe, Japan and North America [6].

The aetiology of KD remains unknown. However, the seasonality and epidemicity seen in areas of high incidence, including Japan, suggest that it is caused by an infectious trigger [7]. The current consensus is that, in some genetically predisposed children, an infectious trigger initiates an abnormal immune response [8,9]. Multiple viral and bacterial pathogens have been suggested as candidates for the trigger, in addition to airborne environmental and fungal triggers [8,10]. Despite the many theories postulated, none have been independently confirmed, and some studies have concluded that KD is likely to be caused by multiple environmental triggers [11].

As the coronavirus disease 2019 (COVID-19) pandemic evolved in early to mid-2020, an increase in cases of children with unusual febrile illnesses, some with features resembling KD, was observed [12]. This new condition, which was later termed “Paediatric Inflammatory Multisystem Syndrome Temporally associated with SARS-CoV-2”, or “Multisystem Inflammatory Syndrome in Children” (PIMS-TS or MIS-C) [12–15], tends to arise several weeks after SARS-CoV-2 infection [14]. The finding of increased KD-like cases after the emergence of a novel viral pathogen raises questions about whether more than one type of trigger might cause KD, and whether KD might represent a constellation of overlapping inflammatory syndromes.

Study of host transcriptomic and proteomic profiles can reveal perturbations caused by infection or inflammation. Comparison of the transcriptional response in different diseases has revealed different host responses to individual pathogens such as TB, bacterial and viral infections [16,17]. Previous studies of host 'omics in the context of KD have characterised the pathways involved in the disease and have established biomarker signatures with diagnostic potential [18,19]. Interrogating multiple 'omic datasets in parallel provides more accurate insights into the molecular dynamics of infection as information captured in one 'omic layer might not necessarily be captured in another 'omic layer.

We explored the host transcriptomic and proteomic profiles of children with KD and those with viral and bacterial infections, aiming to elucidate whether the inflammatory response in KD is more similar to that of a bacterial or viral infection, or indeed neither. We also used the approach to explore the heterogeneity within the transcriptional and translational response of patients with KD.

## 2. Results

### 2.1. Description of Datasets

Whole-blood transcriptomic profiles generated from 414 children were included in the analysis, obtained from children with Kawasaki disease (KD;  $n = 178$ ), confirmed (definite) bacterial infection (DB;  $n = 54$ ), confirmed (definite) viral infection (DV;  $n = 120$ ), and healthy controls (HC;  $n = 62$ ). Two transcriptomic datasets were used. The ‘discovery’ transcriptomic dataset, which was generated by HumanHT-12 version 4.0 BeadChips, was used for all steps of the analysis. The ‘validation’ transcriptomic dataset, which was created by merging two datasets generated by HumanHT-12 version 3.0 and 4.0 BeadChips, was used to test the classifiers trained on the discovery dataset (Table 1).

In addition, proteomic profiles from the plasma or serum of 329 children in the same groups were studied: from children with KD ( $n = 52$ ), DB ( $n = 121$ ) and DV ( $n = 106$ )

infections, and HC ( $n = 50$ ). Liquid chromatography with tandem mass spectrometry (LC–MS/MS) and the SomaScan [20] platform were used to generate the proteomic ‘discovery’ and ‘validation’ datasets, respectively (Table 1). The ‘discovery’ proteomic dataset, generated from plasma samples using LC–MS/MS, was used for all steps of the analysis. The ‘validation’ proteomic dataset, generated from serum samples using the SomaScan platform [20], was used to test the classifiers trained on the discovery dataset.

On both ‘omic levels, the datasets that were used as ‘discovery’ datasets were selected due to their higher number of bacterial and viral samples. There was no overlap between the patients included in the proteomic datasets and those included in the transcriptomic datasets.

KD patients were defined according to American Heart Association (AHA) guidelines [21]. DB patients had a bacterial pathogen identified in a sample from a sterile site. DV patients had a virus identified that was consistent with the presenting syndrome; had no bacteria identified in blood or relevant culture sites; and had C-reactive protein (CRP) levels  $<60$  mg/L. Further details on the clinical definitions used to define the DV and DB groups can be found in the Supplementary Text.

The median ages (months) of KD patients in the transcriptomic discovery and validation datasets were 26 (IQR: 29) and 37 (IQR: 34), respectively. The proportions of male KD patients were 55% and 60% for the transcriptomic discovery and validation datasets, respectively. For the proteomic KD group, the median ages (months) were 30 (IQR: 36) and 16 (IQR: 39) for the discovery and validation datasets, respectively. The proportion of males was 69% for both the discovery and validation datasets (Tables S1 and S2). Table S2 contains clinical information for the KD patients included in the four datasets analysed. The causative pathogens for the patients with bacterial and viral infections from all datasets are shown in Table S3. The median duration of fever when the blood sample was taken for transcriptomic analysis from KD patients was 5 (range of 2–7 days) and 6 days (range of 2–10 days) for the discovery and validation datasets, respectively. For the proteomic KD samples, the median duration of fever when the sample was taken was 7 (range of 3–20 days) and 6.5 days (range of 4–22 days) for the discovery and validation dataset, respectively.

**Table 1.** The datasets used in the analysis. KD, DB, DV and HC are abbreviations for Kawasaki disease, definite bacterial, definite viral, and healthy control, respectively. LC–MS/MS is an abbreviation for liquid chromatography with tandem mass spectrometry. \* = not used in analysis.

Dataset Name	GEO Accession(s)	Platform(s) Used for Generation	KD	DB	DV	HC	Citation(s)
Transcriptomic discovery	GSE73461	Microarray: HumanHT-12 version 4.0	77	31	92	62	[18]
Transcriptomic validation	GSE73462 GSE73463	Microarrays: 1 × HumanHT-12 version 3.0 1 × HumanHT-12 version 4.0	101	23	28	16 *	[19,22]
Proteomic discovery	NA	LC–MS/MS	26	73	75	25	unpublished
Proteomic validation	NA	SomaScan [20]	26	48	31	25	unpublished

## 2.2. Comparison of Kawasaki Disease to Bacterial and Viral Infection

We explored whether the host response during KD is more similar to the host response during bacterial or viral infections using transcriptomic (gene-level) and proteomic data. We first assessed the variance in the discovery datasets using Principal Component Analysis (PCA; Figure S1 and S2). In the transcriptomic dataset, PC1 (29.24%) appears to be capturing lymphocyte number and disease group, with the KD patients located between the bacterial and viral groups. In the proteomic dataset, PC1 (29.18%) appears to be capturing variation caused by age differences, while PC2 (13.39%) and PC3 (10.56%) strongly capture the

disease group effects, with the KD patients grouped together between the clearly separated bacterial and viral groups.

### 2.2.1. Differential Abundance Analysis

Limma [23] was used to identify genes and proteins differentially abundant between each disease group (KD, DB, DV) and healthy controls (HC), whilst accounting for age, sex and, for the transcriptomic dataset, immune cell proportions. Features were considered significantly differentially abundant (SDA) at a FDR of 5%. Differential abundance analysis was applied to 13,035 genes and 344 proteins. For the transcriptomics, 3213, 3124, and 4663 genes were SDA between KD vs. HC, DB vs. HC, and DV vs. HC, respectively. For the proteomics, 113, 125, and 78 proteins were SDA between KD vs. HC, DB vs. HC, and DV vs. HC, respectively. Genes and proteins SDA between KD vs. HC are listed in the Supplementary File S1.

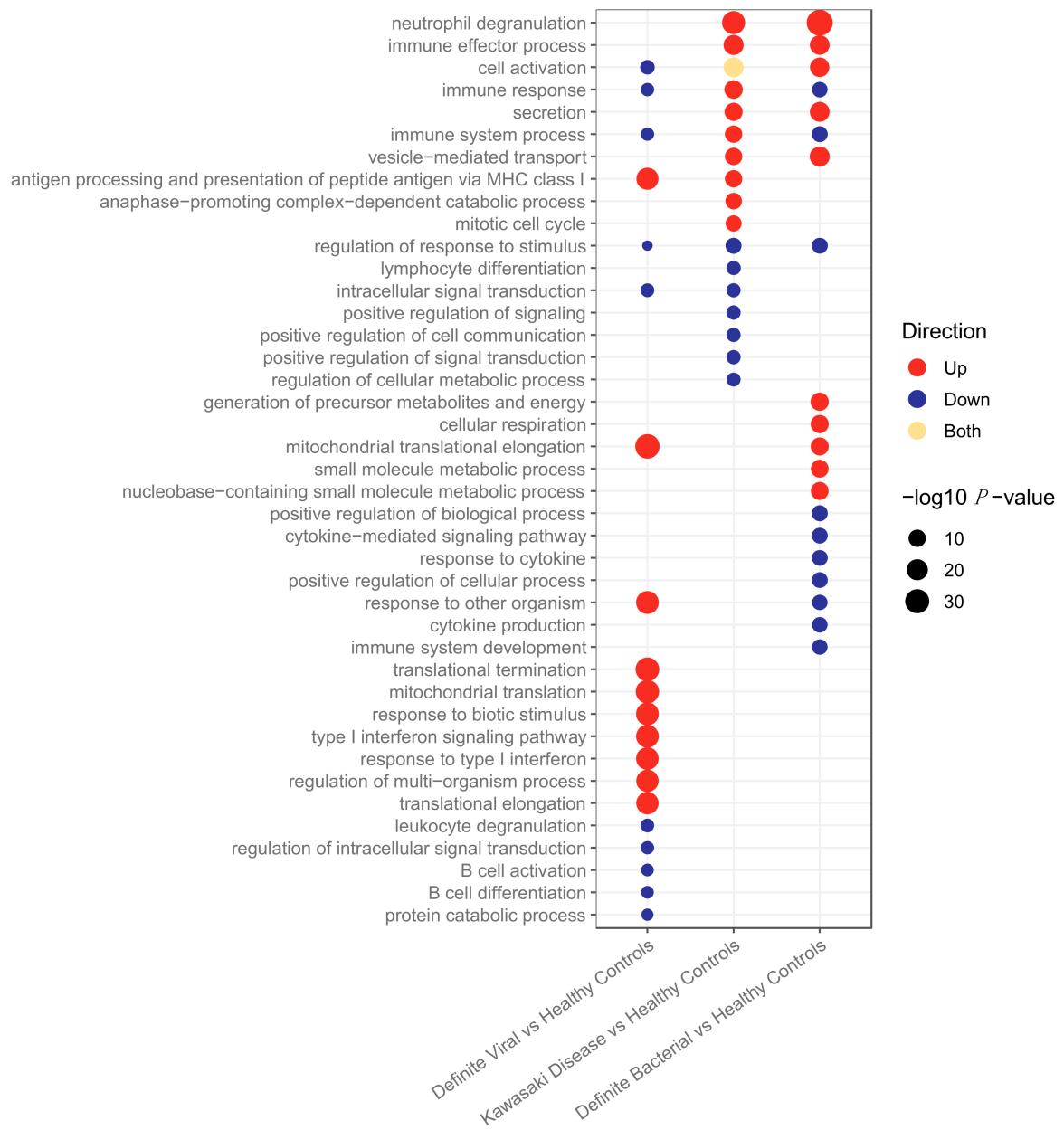
### 2.2.2. Pathway Analysis

The lists of SDA features identified in Section 2.2.1 were subjected to pathway analysis using g:Profiler2 [24] to determine which pathways were upregulated and downregulated in the three disease groups in the discovery datasets for the transcriptomic (Figure 1a) and proteomic (Figure 1b) datasets. The full lists of pathways are provided in Supplementary File S2.

In the transcriptomic pathway analysis, some pathways were found to be enriched across two or three of the disease conditions, whereas others were found in a single condition (Figure 1a). For example, neutrophil degranulation, which was the top pathway in both KD and bacterial infections, and vesicle-mediated transport were both upregulated in KD and bacterial infections, whereas antigen presentation via MHC class I was upregulated in KD and viral infections. Of the top 17 pathways enriched in KD, 6 pathways were also present with concordant directions in the top bacterial pathways and 4 were present with concordant directions in the top viral pathways. Seven were unique to KD.

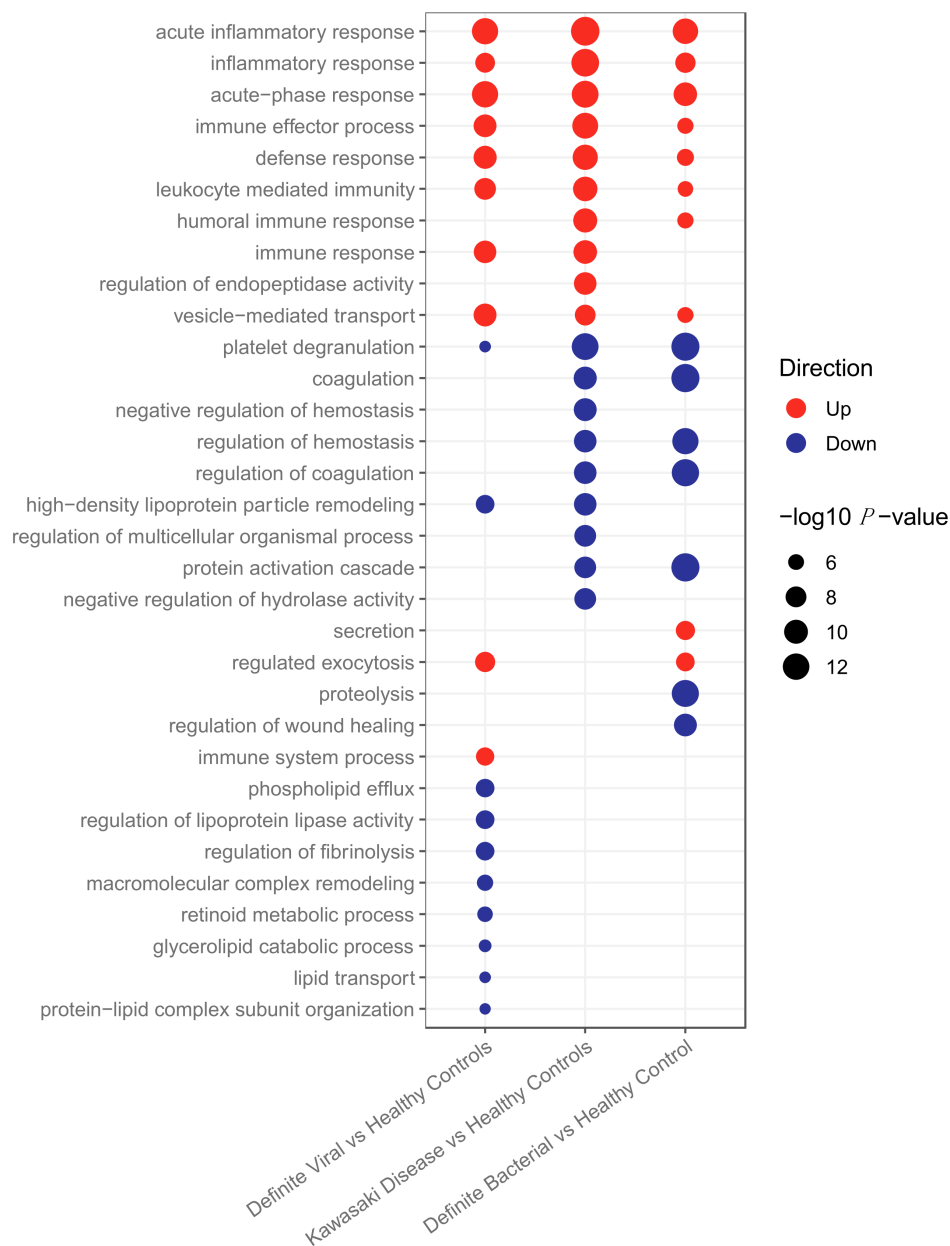
In the proteomics data, all pathways overlapping between KD and either bacterial or viral infections had concordant directions of regulation (Figure 1b). Of the top 19 pathways enriched in KD, 13 were also enriched in bacterial samples, 10 in viral samples, and 4 were unique to KD. Eight of the top 19 KD pathways were enriched in both bacterial and viral samples. All of these are involved in the immune response. The higher frequency of overlapping concordant pathways makes it harder to identify differences between the pathways enriched in the proteomic dataset than the transcriptomic dataset. Overall, there was a much lower number of proteins SDA between KD vs. HC ( $n = 113$ ) than genes SDA between KD vs. HC ( $n = 3213$ ). Furthermore, the total number of proteins remaining following quality control and filtering for missingness ( $n = 344$ ) was much lower than the total number of genes remaining following quality control ( $n = 13,035$ ), which could justify why the differences in pathways enriched between the disease groups are more apparent in the gene expression data.

Three pathways were enriched on both 'omic levels. These were: immune effector process pathway (upregulated in KD and bacterial patients); immune response (upregulated in KD patients); and vesicle-mediated transport (upregulated in KD and bacterial patients).



(a)

Figure 1. Cont.



(b)

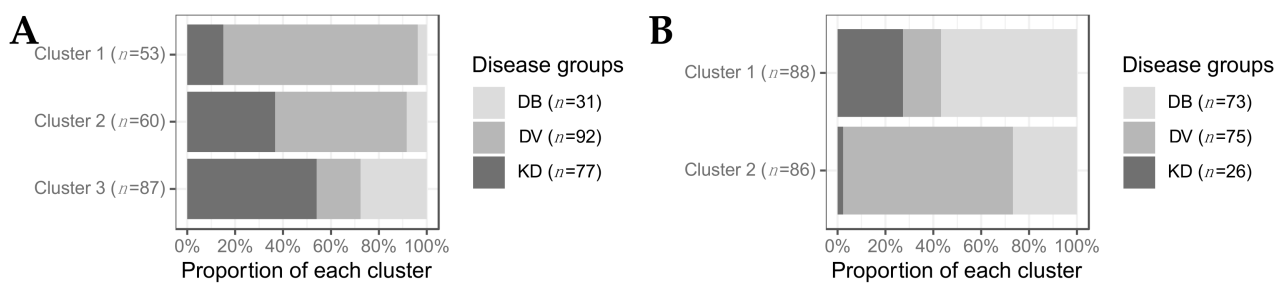
**Figure 1.** (a) Pathways upregulated and downregulated in bacterial, Kawasaki Disease and viral patients compared to healthy controls in the transcriptomic dataset. (b) Pathways upregulated and downregulated in bacterial, Kawasaki Disease and viral patients compared to healthy controls in the proteomic dataset.

### 2.2.3. Clustering

K-Means clustering was used to determine whether the KD patients were more likely to cluster with bacterial or viral patients in the discovery datasets. Prior to clustering analysis, gene expression values were corrected for age, sex and immune cell proportions by taking the residual gene expression values after removing the contributions of these variables. Immune cell proportions were estimated using CIBERSORTx [25], an online tool for estimating immune cell proportions from gene expression data. Without correcting the transcriptomic data for immune cell proportions, clusters formed according to immune

cell proportion (Figures S3 and S4). The same process was performed to remove the contribution of age and sex from the protein abundance values. NbClust [26] was used to determine the optimal number of clusters ( $k$ ). The value of  $k$  most frequently selected across the 12 indices measured by NbClust was selected as the optimal number of clusters for downstream analyses. In the transcriptomic analysis, three clusters were identified as optimal, whereas on the protein level, two clusters were identified as optimal.

We assessed the proportion of KD, bacterial and viral patients in each of the clusters for the transcriptomic (A) and proteomic (B) datasets (Figure 2). In the transcriptomic analysis, an over-representation of viral patients was observed in cluster 1 and, to a lesser extent, cluster 2. An over-representation of bacterial patients was observed in cluster 3 (Figure 2A), resulting in two viral-like clusters and one bacterial-like cluster. Of the 77 transcriptomic KD samples, 47 (61%) belonged to cluster 3, 22 (29%) belonged to cluster 2, and 8 (10%) belonged to cluster 1. In the proteomic analysis, an over-representation of bacterial patients was found in cluster 1, whereas an over-representation of viral patients was observed in cluster 2, leading to one viral-like and one bacterial-like cluster (Figure 2A). Of the 26 proteomic KD samples, 24 (92%) belonged to cluster 1 and 2 (8%) belonged to cluster 2.

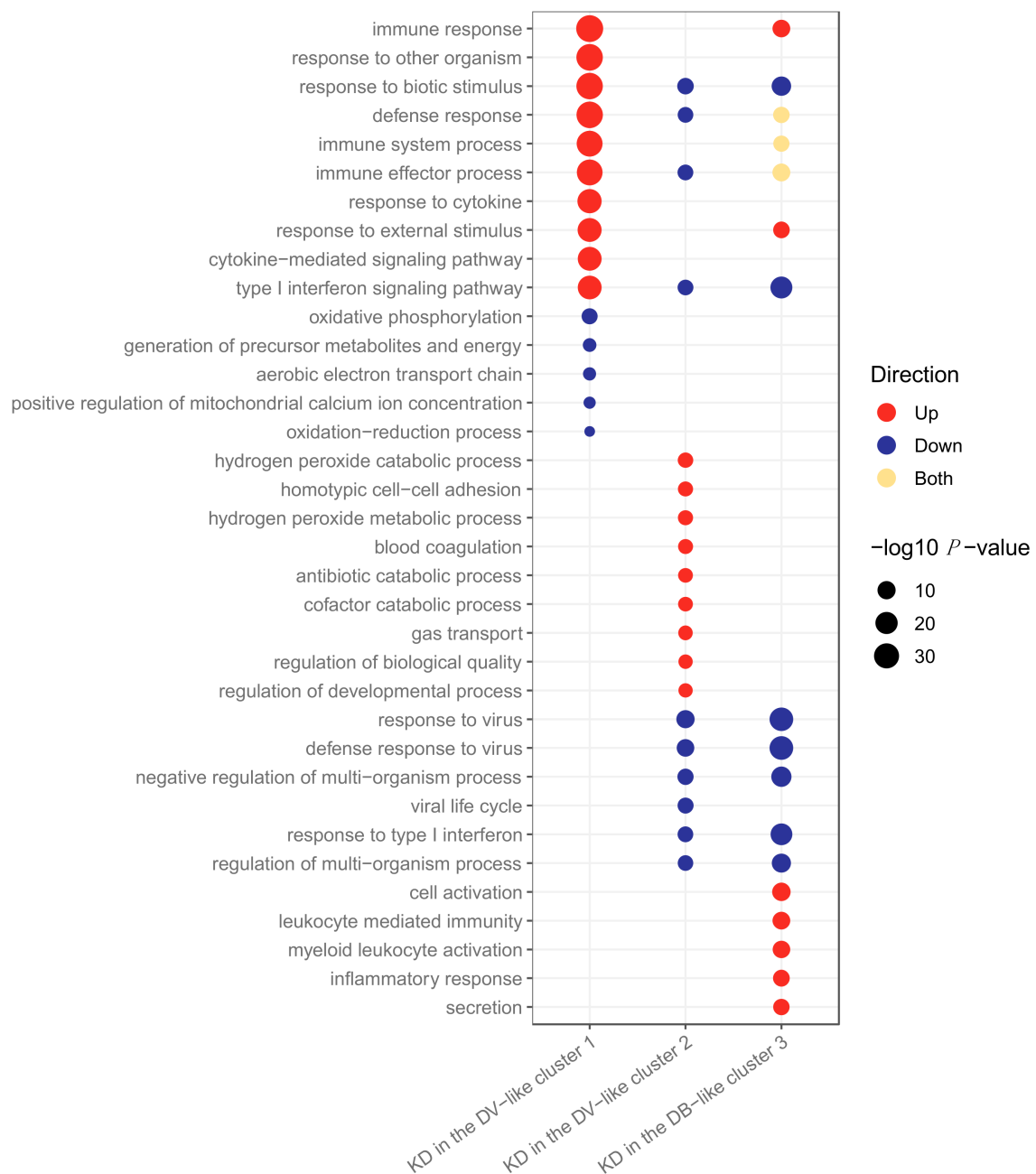


**Figure 2.** The proportion of patients from each disease group in each cluster for transcriptomics (A) and proteomics (B). DB, DV and KD represent definite bacterial, definite viral, and Kawasaki disease.

The association between KD patient cluster membership and various clinical variables was tested. CRP levels ( $p$ -value:  $4.8 \times 10^{-2}$ ) and lymph node swelling ( $p$ -value:  $4 \times 10^{-2}$ ), and peeling ( $p$ -value:  $5 \times 10^{-2}$ ) were significantly associated with cluster membership of KD transcriptomic samples. Higher levels of CRP were found in transcriptomic KD samples in cluster 3 which had the highest proportion of bacterial samples. Out of the 55 patients displaying peeling, 38 were found in cluster 3, as were 17 of the 21 patients with lymph node swelling. No clinical variables were associated with cluster membership in the proteomic dataset. On both 'omic levels, CRP levels were highest in the clusters in which the majority of bacterial samples were found (Figures S5 and S6), as expected since a CRP cut-off of  $<60$  mg/L was required for patients in the DV groups. This pattern was also observed for the WBC counts in the transcriptomic dataset (Figure S5).

Differential abundance analysis was performed to compare feature abundance in the KD samples that fell into different clusters. There were 503 genes SDA between transcriptomic KD samples in cluster 1 vs. clusters 2 and 3, 454 genes SDA between KD samples in cluster 2 vs. clusters 1 and 3, and 651 genes SDA between KD samples in cluster 3 vs. clusters 1 and 2. These lists of SDA genes were subjected to pathway analysis using g:Profiler2 [24] to identify pathways upregulated and downregulated within the clusters (Figure 3). Complete lists of pathways are in Supplementary File S3.





**Figure 3.** Pathways upregulated and downregulated in the KD patients in clusters 1, 2 and 3 for the transcriptomic dataset. Clusters were identified using *K*-Means applied to KD, DB and DV patients. KD, DB and DV represent Kawasaki Disease, definite bacterial, and definite viral, respectively. There were 151, 52 and 137 pathways upregulated in clusters 1, 2 and 3, respectively, and 5, 66 and 137 pathways downregulated in clusters 1, 2 and 3, respectively.

For the transcriptomics, cluster 1 had the highest proportion of viral patients compared to the other clusters (Figure 2A). The majority of the adenovirus (19/23) and influenza (16/23) patients were in cluster 1. Cluster 1 KD patients were characterised by upregulation of anti-viral response pathways, such as interferon and cytokine signalling (Figure 3). In cluster 2, although the majority of patients were viral, their proportion was not quite as high as it was in cluster 1 (Figure 2). The majority of the RSV (15/27) patients were in cluster 2. In the KD patients in cluster 2, various pathways associated with the anti-viral response were downregulated (Figure 3). Cluster 3 had the highest proportion of bacterial patients and KD patients (Figure 2). Similarly to cluster 2, the top pathways downregulated for KD patients in cluster 3 were associated with the anti-viral response, while the inflammatory

response pathway was strongly upregulated, suggesting that the KD patients in this cluster were different to those in cluster 1 and that their response was not as viral-like as those in cluster 1 (Figure 3).

Three pathways—response to biotic stimulus (i.e., a stimulus caused or produced by a living organism), response to other organism and type I interferon signalling—were upregulated in viral transcriptomic samples (Figure 1a) and also in the KD samples in the viral-like cluster 1 (Figure 3). Furthermore, four pathways, including two associated with interferon signalling, were upregulated in viral transcriptomic samples (Figure 1a) and downregulated in the KD samples in clusters 2 and 3 (Figure 3). There were five pathways downregulated in bacterial transcriptomic samples (Figure 1a) and upregulated in KD transcriptomic samples in cluster 1 (Figure 3), including two related to cytokine signalling.

For the proteomic dataset, two proteins were SDA between clusters 1 and 2: serum amyloid A1 (SAA1) and retinol binding protein 4 (RBP4). Both of these proteins have been identified previously as Kawasaki markers, with RBP4 abundance being lower in active KD [27] and SAA1 being elevated in KD [28]. The two KD patients in cluster 2 displayed the opposite pattern, with higher RBP4 and lower SAA1 abundance than the other KD patients.

#### 2.2.4. Classification Using Disease Risk Scores

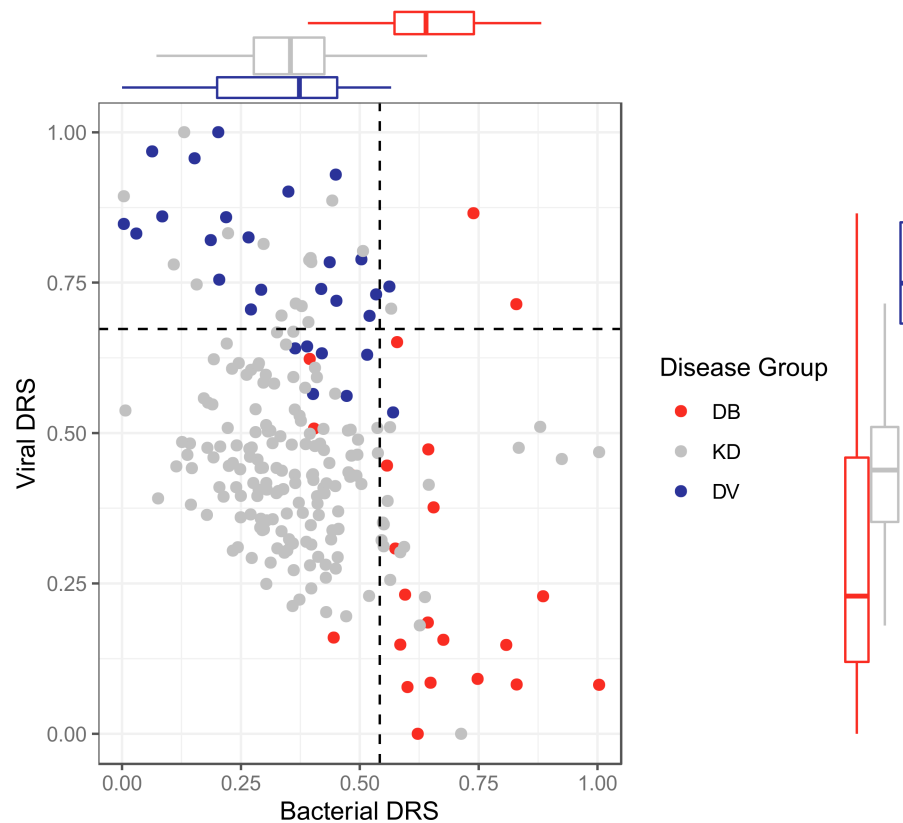
To further assess whether the KD patients elicited more bacterial-like or more viral-like responses, we built two classifiers that returned the probabilities that a patient is bacterial or viral through two separate disease risk scores (DRS). A DRS translates the abundance of features in a discriminatory signature, selected by Lasso [29], into a single value that can be assigned to each individual [16]. Through using two independent classifiers, the possibility of a patient being neither bacterial nor viral was allowed. The classifiers were trained using the 'omic data that was corrected for age, sex and, for the transcriptomic dataset, immune cell proportions.

The Lasso model selected 38 genes for the bacterial classifier, of which 26 had increased abundance and 12 had decreased abundance in bacterial patients compared to viral patients and healthy controls (Table S4). The viral classifier included 32 genes, of which 13 had increased abundance and 19 had decreased abundance in viral patients compared to bacterial patients and healthy controls (Table S5). The classifiers trained in the transcriptomic discovery dataset were tested on bacterial and viral patients from the transcriptomic validation dataset. The bacterial classifier achieved an area under the ROC curve (AUC) of 0.935 (95% CI: 0.869–1) and the viral classifier achieved an AUC of 0.935 (95% CI: 0.856–1).

The Lasso model selected 26 proteins for the bacterial classifier, of which 12 had increased abundance and 14 had decreased abundance in bacterial patients compared to viral patients and healthy controls (Table S6). The viral classifier included 20 proteins, of which 11 had increased abundance and 9 had decreased abundance in viral patients compared to bacterial patients and healthy controls (Table S7). When testing the classifiers trained in the proteomic discovery dataset on bacterial and viral patients from the validation dataset, the bacterial classifier achieved an AUC of 0.925 (95% CI: 0.867–0.984) and the viral classifier achieved an AUC of 0.891 (95% CI: 0.821–0.962). For both 'omic levels, the 90% sensitivity of the classifiers in classifying these samples was used to determine the DRS threshold above which a sample would be classified as bacterial or viral.

The classifiers were applied to KD patients from the discovery and validation datasets for both 'omic levels, resulting in bacterial DRS (DB-DRS) and viral DRS (DV-DRS) for each KD patient (Figures 4 and 5). Classification labels (DB or DV) were assigned to the KD patients using the DB-DRS and DV-DRS thresholds calculated from applying the classifiers to the bacterial and viral patients in the validation datasets (Figure S9). Of the 178 transcriptomic KD samples, 18 (10%) samples had DB-DRS high enough to be classified as bacterial and 16 (9%) samples had DV-DRS high enough to be classified as viral. 145 (81%) samples did not achieve DB-DRS nor DV-DRS sufficiently high to lead

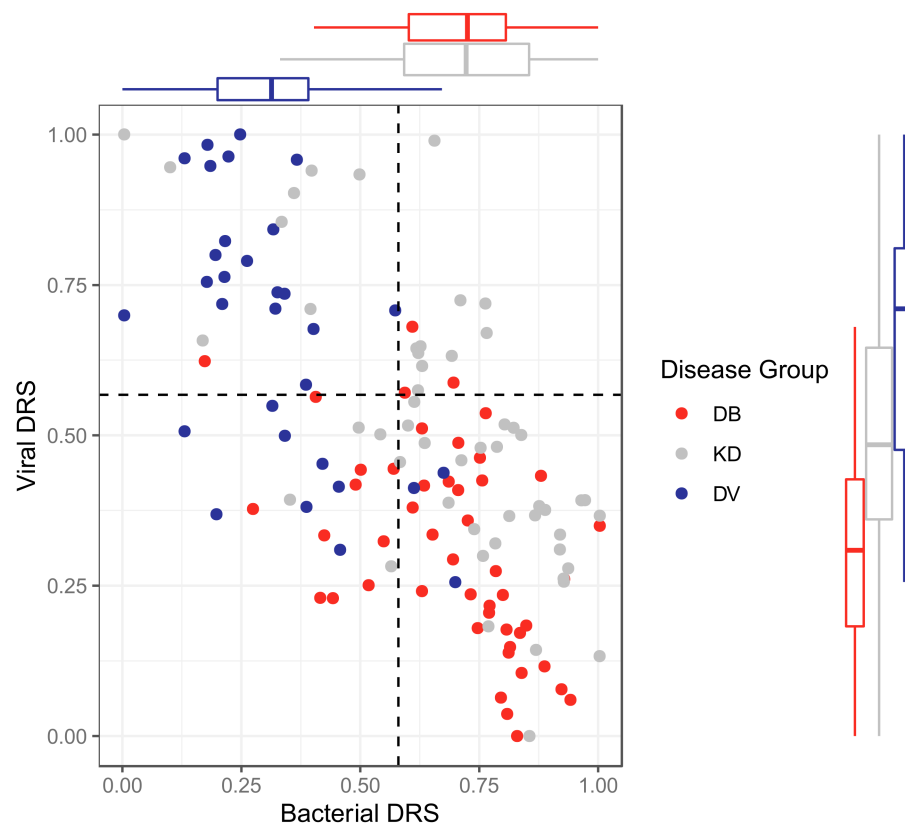
to bacterial or viral classification, and 1 sample was classified as both bacterial and viral (Figure 4). Of the 52 proteomic KD samples, 40 (78%) achieved DB-DRS high enough to be classified as bacterial and 18 (35%) achieved DV-DRS high enough to be classified as viral. 10 (19%) proteomic KD samples achieved DB-DRS and DV-DRS high enough for them to be classified as both bacterial and viral, and 4 (7.7%) were classified as neither (Figure 5).



**Figure 4.** Bacterial DRS (DB-DRS) plotted against viral DRS (DV-DRS) for KD (discovery and validation), DB (from the validation cohort) and DV (from the validation cohort) patients from the transcriptomic datasets. Boxplots are shown for each disease group. KD, DB and DV represent Kawasaki Disease, definite bacterial, and definite viral, respectively.

To further examine the ‘omic profiles of the KD patients with DB-DRS and DV-DRS too low for them to be classified as either bacterial or viral, we performed pathway analysis on the genes or proteins SDA between these KD patients and healthy controls. Amongst the pathways upregulated on the transcriptomic level, were ‘defence response to fungus’ ( $p$ -value:  $7 \times 10^{-08}$ ) and ‘response to fungus’ ( $p$ -value:  $7 \times 10^{-07}$ ).

The associations of DB-DRS, DV-DRS, bacterial classification as predicted from the DB-DRS, and viral classification as predicted from the DV-DRS, with various clinical variables were tested for KD samples from both ‘omic levels. In the transcriptomic KD samples, clinical measurements of CRP were positively associated with DB-DRS ( $p$ -value: 0.002) and bacterial classification ( $p$ -value: 0.0001), and negatively associated with DV-DRS ( $p$ -value: 0.002) and viral classification ( $p$ -value: 0.023). In the proteomic KD samples, CRP levels were significantly positively associated with DB-DRS ( $p$ -value: 0.013) and bacterial classification ( $p$ -value: 0.007). Peeling was significantly associated with higher DB-DRS on both ‘omic levels (transcriptomic  $p$ -value: 0.041, proteomic  $p$ -value: 0.007). Strawberry tongue was significantly associated with a low score on the transcriptomic DV-DRS ( $p$ -value: 0.045).



**Figure 5.** Bacterial DRS (DB-DRS) plotted against viral DRS (DV-DRS) for KD (discovery and validation), DB (from the validation cohort) and DV (from the validation cohort) patients from the proteomic datasets. Boxplots are shown for each disease group. KD, DB and DV represent Kawasaki Disease, definite bacterial, and definite viral, respectively.

For the KD patients from the discovery datasets, the associations between DB-DRS or DV-DRS and the cluster membership of patients were tested. Transcriptomic KD sample cluster membership was significantly associated with DB-DRS ( $p$ -value: 0.005) and DV-DRS ( $p$ -value: 0.0006), with a stepwise increase in DB-DRS and decrease in DV-DRS from clusters 1 to 3, where cluster 1 was the most viral-like cluster, and cluster 3 was the most bacterial-like cluster. Proteomic KD sample cluster membership was significantly associated with DB-DRS ( $p$ -value: 0.002) and DV-DRS ( $p$ -value: 0.023), with higher DB-DRS and lower DV-DRS in KD patients in cluster 1, where cluster 1 was the more bacterial-like cluster and cluster 2 was the more viral-like cluster.

### 2.3. Clustering of Kawasaki Disease Patients Alone

We performed unsupervised clustering for the KD patients from the discovery datasets to explore the natural patient stratification formed in the absence of bacterial and viral comparator patients. For both 'omic levels, 3 clusters were optimal, as determined by NbClust [26]. The clusters were identified using the 'omic data that was corrected for age, sex and, for the transcriptomic dataset, immune cell proportions. Of the 77 transcriptomic KD samples, 32 (41%) were in cluster 1 (cluster KD1-T), 23 (30%) were in cluster 2 (cluster KD2-T), and 22 (29%) were in cluster 3 (cluster KD3-T). Of the 26 proteomic KD samples, 4 (15%) were in cluster 1 (cluster KD1-P), 7 (27%) were in cluster 2 (cluster KD2-P), and 15 (58%) were in cluster 3 (cluster KD3-P).

There was high overlap between the samples in cluster KD1-T and those in the transcriptomic bacterial-like cluster 3 described in Section 2.2.3 (Figure S10). All except one of the samples found previously in the transcriptomic viral-like cluster 1 were found in cluster KD2-T. The majority ( $n = 14$ ; 64%) of the samples in KD3-T were also found in

transcriptomic cluster 2. On the proteome level, in Section 2.2.3, all KD samples except two clustered together in cluster 1. However, the two remaining samples that were previously in cluster 2 were not assigned to the same cluster.

The association between cluster membership and various clinical variables was tested. CRP levels were significantly associated with cluster membership for both 'omic layers (transcriptomics  $p$ -value: 0.041, proteomics  $p$ -value: 0.010). Furthermore, coronary artery aneurysm (CAA) formation was significantly associated with cluster membership in the proteomic dataset ( $p$ -value: 0.020) with 13 of the 21 patients known to not have CAAs being in cluster KD3-P. On the transcriptomic level, the highest WBC counts and CRP levels were in cluster KD1-T, and on the proteomic level, WBC counts and CRP levels were highest in clusters KD2-P and KD1-P, respectively (Figure S7 and Figure S8).

The associations between DB-DRS or DV-DRS and cluster membership of KD patients when clustered alone was tested. The transcriptomic KD samples' cluster membership was significantly associated with DB-DRS ( $p$ -value: 0.006) with the highest DB-DRS in cluster KD1-T. Although the association between transcriptomic KD samples' cluster membership and DV-DRS was not significant, the highest DV-DRS values were observed in cluster KD2-T. There were no significant associations between the proteomic KD samples' cluster membership and DB-DRS or DV-DRS.

Differential abundance analysis was performed on the patients that fell into different clusters. For the transcriptomics, there were 494 genes SDA between cluster KD1-T vs. clusters KD2-T and KD3-T, 461 genes SDA between cluster KD2-T vs. clusters KD1-T and KD3-T, and 320 genes SDA between cluster KD3-T vs. clusters KD1-T and KD2-T. For the proteomics, 42 proteins were SDA between cluster KD1-P vs. clusters KD2-P and KD3-P, 25 proteins were SDA between cluster KD2-P vs. clusters KD1-P and KD3-P, and 38 proteins were SDA between cluster KD3-P vs. clusters KD1-P and KD2-P. These lists of SDA features were subjected to pathway analysis using g:Profiler2 [24] to identify pathways upregulated and downregulated within the clusters (Figure 6). Complete lists of pathways are found in Supplementary Files S4–S5.

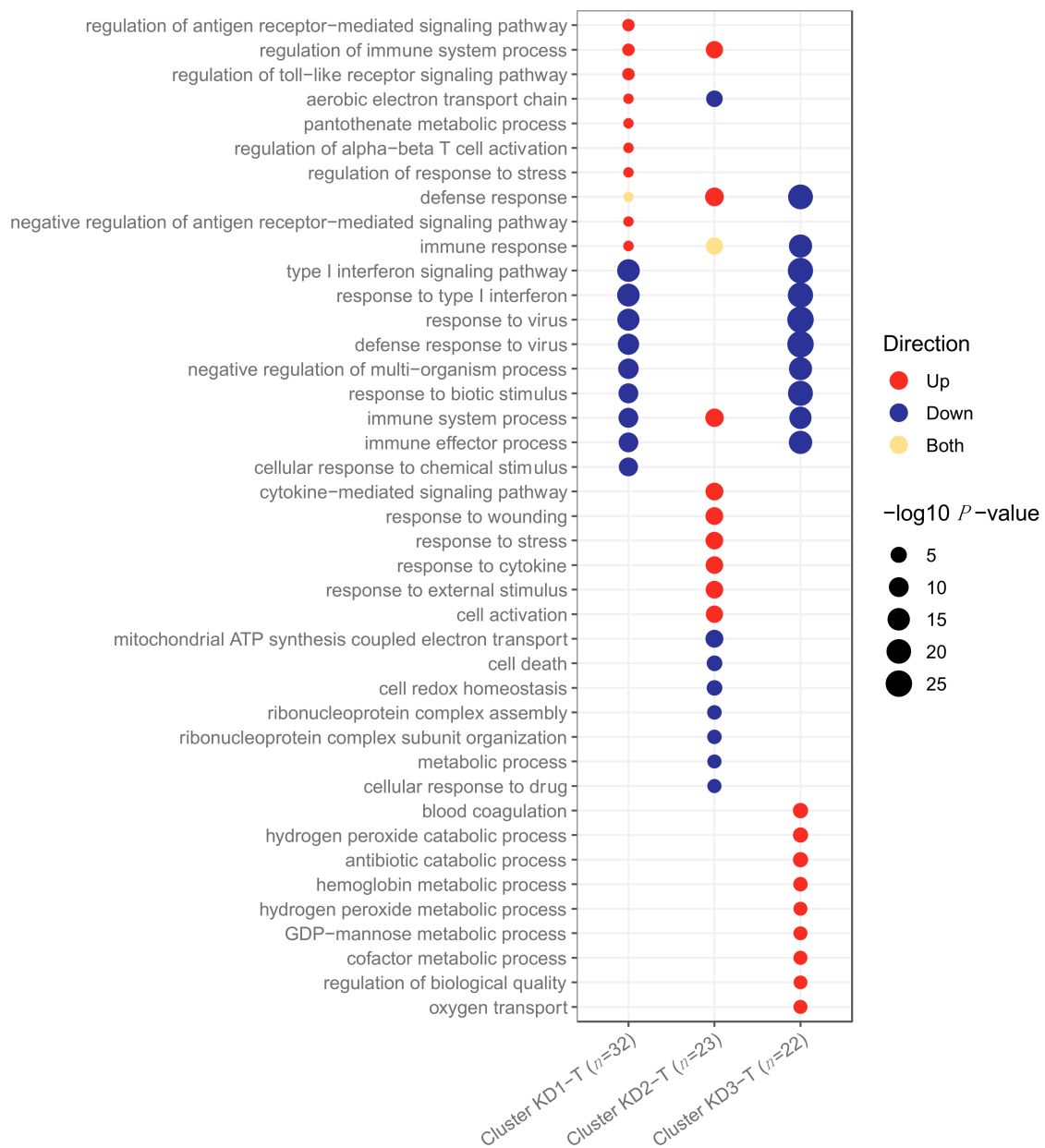
In the transcriptomic analysis (Figure 6a), cluster KD2-T had features in common with an anti-viral response, whilst the others did not. Many pathways associated with the anti-viral response were downregulated in clusters KD1-T and KD3-T, whilst patients in cluster KD2-T were characterised by the upregulation of viral pathways, including those associated with cytokine signalling.

The response to biotic stimulus and type I interferon signalling pathways were previously identified as being upregulated in viral transcriptomic samples (Figure 1a) and in KD samples in the viral-like cluster 1 when K-Means was applied to KD, DB and DV (Figure 3). These pathways were downregulated in clusters KD1-T and KD3-T (Figure 6a), indicating that the transcriptomic response in these samples was less viral-like than samples in cluster KD2-T.

Four pathways previously identified as being downregulated in bacterial transcriptomic samples (Figure 1a), including two pathways associated with cytokine signalling, were upregulated in cluster KD2-T. In addition, six pathways upregulated in cluster KD2-T had already been identified as being upregulated in the viral-like cluster 1 identified previously (Figure 3). Five pathways of the nine top upregulated pathways in cluster KD3-T (Figure 6a) were also upregulated in cluster 2 when K-Means was applied to KD, DB and DV (Figure 3). These were blood coagulation, hydrogen peroxide catabolic process, antibiotic catabolic processes, hydrogen peroxide metabolic process and regulation of biological quality.

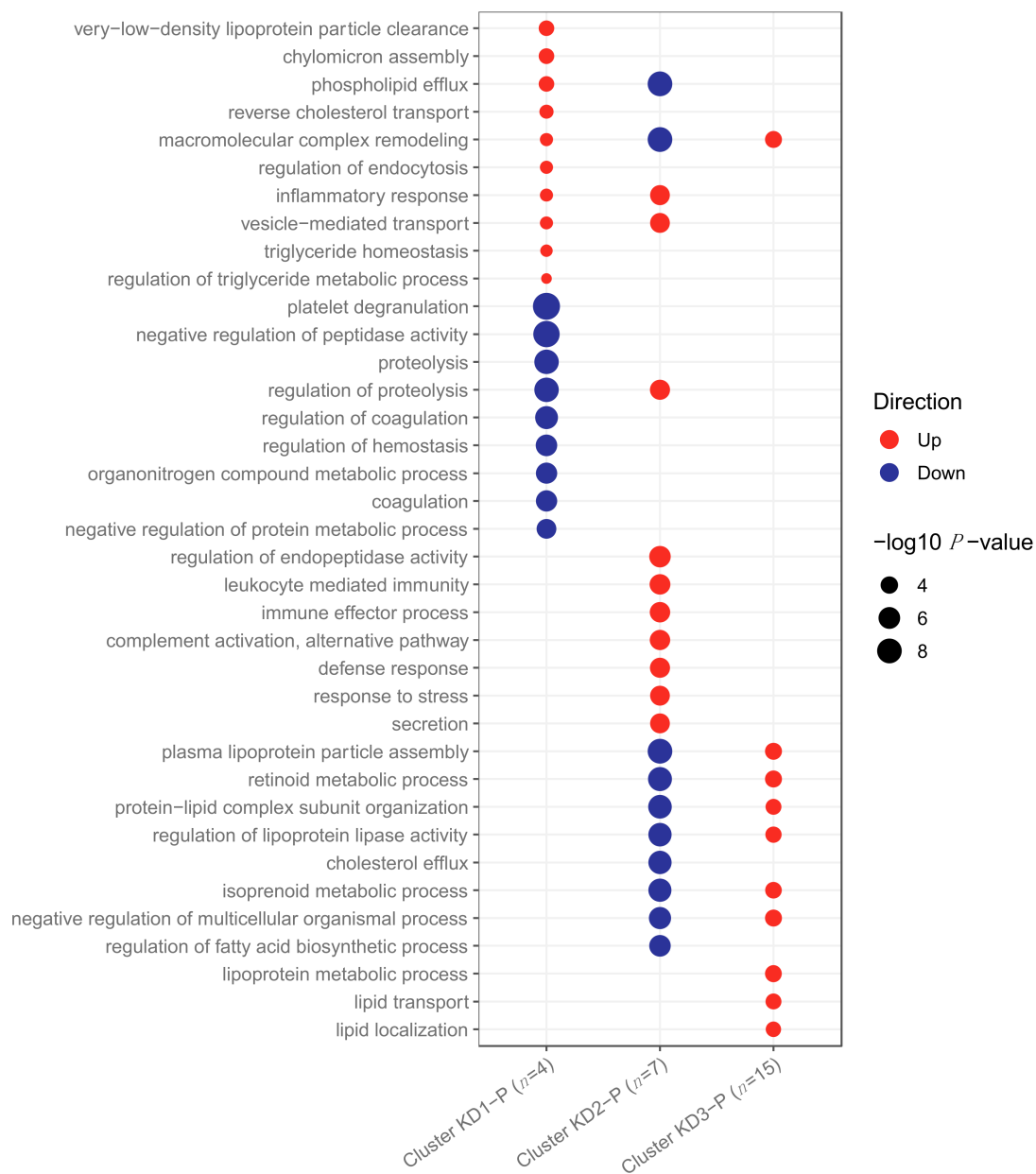
In the proteomic analysis (Figure 6b), one of the KD clusters had features in common with the anti-viral response, whilst another KD cluster was more bacterial-like. Amongst the top pathways enriched in cluster KD1-P and KD2-P were pathways involved in inflammation. The top pathways enriched in cluster KD3-P were associated with lipids. Of the 37 pathways enriched in the proteomic KD samples (Figure 6b), 21 were previously identified as enriched in proteomic samples (Figure 1b). Of these, six were enriched in

proteomic viral samples (Figure 1b) and samples in cluster KD2-P (Figure 6b) with concordant directions. Furthermore, seven pathways enriched in cluster KD1-P (Figure 6b) were also enriched in bacterial proteomic samples (Figure 1b) with concordant directions. These results suggest that cluster KD1-P is a more bacterial-like cluster, whereas cluster KD2-P is a more viral-like cluster. Some pathways were enriched on both 'omic levels, including those associated with blood coagulation, the response to stress and immune effector processes.



(a)

Figure 6. Cont.



(b)

**Figure 6.** (a) Pathways upregulated and downregulated in transcriptomic KD patients between clusters. Clusters were identified by K-Means ran on KD patients alone. There were 24, 118 and 24 pathways upregulated in clusters KD1-T, KD2-T and KD3-T, respectively, and 94, 68 and 75 pathways downregulated in clusters KD1-T, KD2-T and KD3-T, respectively. KD represents Kawasaki Disease. (b) Pathways upregulated and downregulated in proteomic KD patients between clusters. Clusters were identified by K-Means ran on KD patients alone. There were 77, 94 and 61 pathways upregulated in clusters KD1-P, KD2-P and KD3-P, respectively, and 64, 104 and 53 pathways downregulated in clusters KD1-P, KD2-P and KD3-P, respectively. KD represents Kawasaki Disease.

### 3. Discussion

Although the cause of Kawasaki disease has not been identified, there is growing clinical, epidemiological, and immunological evidence that it may be caused by different infectious triggers, with data pointing to bacteria, viruses or fungi. We explored the transcriptomes and proteomes of children with KD and definite bacterial and viral infections,

using multiple approaches to compare the host response to these diseases to the response during KD. We found that there was a diversity of responses in the proteomic and transcriptomic profiles of KD patients, suggesting that KD is not a homogenous condition, and that whilst some patients had a more viral- or bacterial-like profile, the majority were defined as neither bacterial nor viral when their transcriptomic response was mapped onto viral and bacterial disease risk scores (DRS).

Within the host response profiles, some elements of KD appeared more viral-like and some elements appeared more bacterial-like. This is shown through overlapping pathways that were enriched in KD and either bacterial or viral infections. For example, the antigen presentation via MHC class I pathway was upregulated in KD and viral infections on the transcriptomic level. Major histocompatibility complex (MHC) molecules are expressed on the cell surface to present antigenic peptides to T cells, and their expression is increased by a broad range of immune activators including interferons [30,31]. The finding of upregulated MHC class I expression in KD and viral patients may reflect interferon-induced activation in these groups. Additionally, on the transcriptomic level, KD and bacterial infections share neutrophil degranulation as their most upregulated pathway. Neutrophils are the first responders to infection and inflammation, and the expansion and activation of the neutrophil population is a characteristic feature of acute KD. In the initial days of KD illness, there is an intense inflammatory response with neutrophil leucocytosis [32]. Studies have found elevated levels of human neutrophil elastase and IL-8, a C-X-C chemokine that activates neutrophils [33,34].

The host response during KD is highly heterogenous, as demonstrated through the enrichment of certain pathways in the KD patients in different clusters when *K*-Means was applied to KD, bacterial and viral transcriptomic samples. For example, anti-viral response pathways were upregulated in KD patients in the majority viral cluster 1 and downregulated in KD patients in clusters with decreasing numbers of viral samples (cluster 2, 3), relative to each other. In the majority bacterial cluster 3, pathways associated with the inflammatory response were upregulated. The heterogeneity of the host response during KD was also apparent when *K*-Means was applied to KD patients alone. Three distinct clusters were identified on both 'omic levels, and, in each cluster, a distinct set of pathways was enriched. The range of pathways enriched in the different clusters further demonstrate the heterogeneity in the host response during KD, with some clusters enriched for viral response pathways and some clusters enriched for bacterial response pathways. Unsurprisingly, amongst the patients clustering in the more bacterial-like clusters, their DB-DRS tended to be higher, and amongst the patients clustering in the more viral-like clusters, their DV-DRS tended to be higher.

The two different approaches to clustering (with and without bacterial and viral comparator samples) produced similar clusters of KD patients, providing reassurance that the clusters described here are biologically meaningful. Despite the similarities, however, the clusters identified in Sections 2.2.3 and 2.3 were not completely identical, indicating that the inclusion of well-characterised bacterial and viral patients adds further insights to the solely data-driven KD-based analysis.

Although there are shared features between the response to KD and both bacterial and viral infections, the distinct pathways enriched in each disease group demonstrate the variation in the molecular host response; the distinctiveness of the responses is also supported by the ability of RNA and protein signatures to discriminate KD from bacterial and viral infections [18,19,35]. These differences between the response during KD and the responses to bacterial and viral infections suggest that KD may be triggered by a novel process not typical of either common bacterial or viral infections. Despite the host 'omics profiles' heterogeneity observed in KD, commonalities are also shown.

A two-way classifier approach highlighted that it is not a simple dichotomous question as to whether the response during KD more closely resembles the responses to bacterial or viral infections, when focusing on key discriminatory features. We found that 145 of the 178 transcriptomic KD samples were not assigned DRS high enough for them to be



classified as either bacterial or viral, and amongst pathways upregulated in these KD patients compared to healthy controls were two pathways associated with the fungal response. This finding is intriguing, given the evidence suggesting that KD could be caused by a fungal trigger that has been reported elsewhere [36,37].

The heterogeneity and the different clusters of responses to KD which have elements shared with bacterial, viral or fungal responses, could indicate multiple microbial triggers of KD, as has been suggested by Rypdal et al. [11]. An alternative explanation for the heterogeneity observed here in the response during KD could be that a single pathogen that causes KD leads to heterogeneous responses in different hosts, as has been observed in children infected with SARS-CoV-2, where many children remain asymptomatic, some experience severe inflammation [38,39], and some develop PIMS-TS/MIS-C [13,14,40]. Variations in the host condition, such as epigenetic differences and differences in prior pathogen exposure, could cause the spectrum of host responses to KD observed here. Differences in host genetics could also be responsible for the heterogeneity in host response during KD as the severity of KD, including the formation of CAAs, is already known to be impacted by the host's genetic background [41].

This study has certain limitations. The proteomic discovery dataset was a lower-dimensional dataset ( $n = 867$ ) than the transcriptomic discovery dataset ( $n = 47,323$ ) with high rates of missingness, as is common in quantitative proteomics. Only proteins with no missingness were used for the clustering and classification, so key proteins for distinguishing KD could be absent from the analysis. On the proteomic level, many pathways were enriched in multiple disease groups (Figure 1b), making it difficult to identify a disease-specific pathway signature. This could be caused by plasma samples, which were used in this dataset, capturing a noisy signal due to the release of substances from various tissues into the bloodstream. The proteomic response during KD shared more similarities with the proteomic response to bacterial infection, with more pathways overlapping between KD and bacterial infections (Figure 1b) and all but two KD proteomic samples clustering with bacterial proteomic samples (Figure 2). This follows observations of striking clinical similarities between KD and bacterial streptococcal and staphylococcal toxic shock syndromes [42,43], and could reflect the hypothesis that the proteome is closer to the observed phenotype than the transcriptome [44].

Although bacterial patients with known viral coinfections have been removed from the analysis, it is impossible to say with confidence that an individual does not have a coinfection, the presence of which could falsely increase heterogeneity in the host response in a given disease group. Coinfection is common in KD, with one study identifying confirmed infections in a third of KD patients [45]. Despite being unable to rule-out that some KD patients included had co-incident viral or bacterial infections, we found that most KD transcriptomic samples were neither classified as viral nor bacterial when the respective DRS scores were applied. Amongst the patients classified as bacterial or viral, it is possible that some patients could be suffering from an intercurrent infection in addition to KD.

There are variations in the range of bacterial and viral pathogens and the severity of illness represented in the two 'omic datasets. The bacterial and viral patients included in the transcriptomic datasets and the proteomic validation dataset were more severely unwell than those included in the proteomic discovery dataset (Table S3) due to the inclusion criteria of the studies to which they were recruited. The KD patients included in the transcriptomic dataset were collected from San Diego, CA, USA, whereas the KD patients included in the proteomic dataset were collected from London, UK, although the same case definition was used. There remains no diagnostic test for KD, thus some KD patients presented here may have unrecognised alternative diagnoses. The two KD samples in the proteomic dataset that cluster separately (Figure 2B) and are distinguished from the other KD samples by their levels of SAA1 and RBP4, two previously identified KD markers [27,28], are possible examples of this.

## 4. Materials and Methods

### 4.1. Patient Recruitment

All samples were obtained from patients with written parental informed consent. Case definitions can be found in the Supplementary Text. The definite bacterial (DB), definite viral (DV), healthy control (HC) and Kawasaki disease (KD) samples used in the transcriptomic discovery and validation datasets were recruited in the United Kingdom and Spain as part of the IRIS (Immunopathology of Respiratory, Inflammatory and Infectious Disease; NIHR ID 8209) and GENDRES (Genetic, Vitamin D, and Respiratory Infections Research Network; [gendres.org](http://gendres.org)) studies [17,22] and in the United States through the US-Based Kawasaki Disease Research Center Program ([medschool.ucsd.edu/som/pediatrics/research/centers/kawasaki-disease/pages/default.aspx](http://medschool.ucsd.edu/som/pediatrics/research/centers/kawasaki-disease/pages/default.aspx)).

The DB, DV and HC samples used in the proteomic discovery and validation datasets were enrolled in the EUCLIDS (European Union Childhood Life-Threatening Infectious Disease Study; 11/LO/1982) study [46] and the PERFORM (Personalised Risk assessment in Febrile illness to Optimise Real-life Management across the European Union) study ([perform2020.org/](http://perform2020.org/); 16/LO/1684). KD samples used in the proteomic datasets were recruited from the ongoing UK Kawasaki study “Genetic determinants of Kawasaki Disease for susceptibility and outcome” (13/LO/0026). This study recruits acutely unwell children with KD during hospital admission in participating hospitals around the UK.

### 4.2. Data Generation

#### 4.2.1. Transcriptomic Datasets

The transcriptomic discovery dataset was generated from whole-blood samples obtained from KD patients, healthy controls, and patients with bacterial and viral infections using the HumanHT-12 version 4.0 (Illumina, San Diego, CA, USA) microarray [18]. In order to obtain a transcriptomic validation dataset containing the same disease groups as the transcriptomic discovery dataset, two datasets were merged. One dataset contained gene expression values (HumanHT-12 version 4.0 [Illumina, San Diego, CA, USA] microarray) from whole-blood samples obtained from acute and convalescent KD samples [19]. The other dataset consisted of gene expression values (HumanHT-12 version 3.0 [Illumina, San Diego, CA, USA] microarray) from whole-blood samples obtained from patients with bacterial and viral infections [22]. For all three independent microarray experiments, one batch of samples was processed, and samples were randomly positioned across the arrays.

#### 4.2.2. Proteomic Datasets

The proteomic discovery dataset was generated from plasma samples using LC-MS/MS. Full details of the experimental protocol are in the Supplementary Text. The proteomic validation dataset was generated from serum samples using the SomaScan (SomaLogic, Boulder, CO, USA) aptamer-based platform [20]. Prior to pre-processing, 867 proteins were measured in the discovery dataset (LC-MS/MS) and 1300 in the validation dataset (SomaScan). Samples in the proteomic validation dataset were split across three plates with KD, DB, DV and HC samples present on each plate in relative proportions.

### 4.3. Statistical Methods

All analysis was conducted using the statistical software R (R version 3.6.1, [47]). Code used for the analytical pipeline described here is found at [github.com/heather-jackson/KawasakiDisease\\_IJMS](https://github.com/heather-jackson/KawasakiDisease_IJMS). Note, the code is signposted for the transcriptomic datasets but can be modified for other ‘omic levels.

#### 4.3.1. Pre-Processing of Gene Expression Data

Background correction, robust spline normalisation (RSN), and log<sub>2</sub> transformation were applied to the raw discovery gene expression dataset using the R package lumi [48]. Probes were retained if at least 80% of samples in each comparator group had a detection *p*-value < 0.01. Low variance probes and those significantly associated with the UCSD

recruitment site were removed. Bacterial samples with known viral coinfections were removed from the analysis at this stage to ensure that the signal from the bacterial samples was not diluted. KD samples that had been administered IVIG treatment were also removed at this stage, but their inclusion was irrespective of coincident viral or bacterial detection, for which data was not available. A KD sample previously identified as an outlier [18] was removed.

As mentioned, two microarray gene expression datasets were merged to form the validation dataset. Background subtraction and RSN normalisation were applied to these two datasets independently, using the R package lumi [48], prior to using ComBat [49] to remove the batch effects in the merged dataset [18].

#### 4.3.2. Pre-Processing of Protein Abundance Data

The raw discovery dataset files generated by LC–MS/MS were processed using MaxQuant (1.6.10.43) [50]. With matching between runs activated. Relative quantification was performed using the MaxLFQ algorithm [51]. The resulting LFQ values were log<sub>2</sub>-transformed. Bacterial samples with known viral coinfections were removed from the analysis at this stage to ensure that the signal from the bacterial samples was not diluted. Protein groups were removed if they were identified as contaminants, or if they were missing in over 90% of samples in each disease group.

The proteomic validation dataset was generated from the SomaScan platform [20]. Quality control steps used scale factors returned from the SomaScan platform to correct for variations in aptamer hybridisation efficiency, inter- and intra-assay variability, variability in the starting quantities of proteins, and plate effects. Further batch effect corrections were carried out using COCONUT normalisation [52].

#### 4.3.3. Comparison of Kawasaki Disease to Bacterial and Viral Infections Differential Abundance Analysis

Differential abundance analysis was carried out to compare the overall transcriptomic and proteomic responses to KD, definite bacterial (DB) and definite viral (DV) infections. The degree to which genes and proteins were differentially abundant between KD and healthy controls (HC), DB and HC, and DV and HC was quantified using Limma [23] on the transcriptomic and proteomic discovery datasets separately. Age and sex were included as covariates for both datasets. Immune cell proportions, calculated using the online CIBERSORTx portal [25], were used as additional covariates for the transcriptomic dataset. The immune cell proportions included were lymphocytes, neutrophils, monocytes, mast cells and eosinophils. Features were considered significantly differentially abundant (SDA) at a false discovery rate (FDR) [53] of 5%.

#### Pathway Analysis

The pathways upregulated and downregulated in KD, DB, and DV samples were identified from the lists of SDA features identified for each disease group in as outlined in Section Differential Abundance Analysis. Pathways were identified using g:Profiler2 [24] and redundancy in the pathways identified was removed using REVIGO [54].

#### Clustering Analysis

K-Means clustering [55] was applied separately to transcriptomic and proteomic discovery datasets. Healthy controls were excluded as we were only interested in the clustering of KD with pathological patients. For the proteomic dataset, only proteins with no missing data points were used ( $n = 106$ ).

To explore the effects of sex and age on clustering in the proteomic dataset, the contribution of these variables was removed by regressing out their effects on every protein and taking the residual values as the 'corrected' abundance. This process was also followed in the transcriptomic dataset, but the contributions of the immune cell proportions listed in Section Differential Abundance Analysis were also removed. Prior to clustering, and

after correction, features were removed if their variance was lower than 0.25. To determine the optimal number of clusters ( $k$ ) for each corrected and non-corrected dataset, the R package NbClust [26] was used, with 12 indices tested. The indices tested were: KL [56], CH [57], Hartigan [58], McClain [59], Dunn [60], SDIndex [61], SDbw [62], C-Index [63], Silhouette [64], Ball [65], Ptbiserial [66,67] and Ratkowsky [68]. The number of clusters tested by NbClust ranged between 2 and 10 clusters. The most frequently selected  $k$  by the 12 indices was used for downstream analyses. The lowest  $k$  selected the most frequently was taken in cases where there were multiple values of  $k$  selected the most frequently.

Once clusters were identified, features that were SDA (5% FDR) between KD samples in the different clusters were identified. Pathway analysis was done using these lists of SDA features to determine the pathways upregulated and downregulated in KD samples in the different clusters. The R package g:Profiler2 [24] was used for pathway analysis, with pathways with  $p$ -values  $< 0.01$  considered significant. Redundancy in the pathways identified was removed using REVIGO [54].

The associations between cluster membership and various clinical variables were tested. For categorical variables, Fisher's Exact test was used. For continuous variables, one-way ANOVA was used.  $p$ -values  $< 0.05$  were considered significant. The categorical variables tested in both datasets were: strawberry tongue (yes/no/unknown); lymph node swelling (yes/no/unknown); and peeling (yes/no). Continuous variables tested in both datasets were: levels of C-reactive protein (CRP); month of year; and the duration of fever at sampling. Coronary artery aneurysms (CAA) information was available only as a dichotomous variable for the patients submitted for proteomic analysis. For the patients submitted for transcriptomic analysis, maximal coronary artery Z-scores were available and were used instead.

### Classification

Two independent classifiers were built for each 'omic dataset. One classifier was for classifying DB patients (DB classifier), and the other was for classifying DV patients (DV classifier). The DB classifiers and DV classifiers were trained on features SDA between DB vs. DV and HC patients combined, and DV vs. DB and HC patients combined, respectively. Only features present in both datasets (discovery and validation) were used for training the classifiers. The discovery datasets that were corrected for age, sex, and for transcriptomics, immune cell proportions, were used to train the classifiers. The validation datasets were also corrected for age, sex and, for the transcriptomic validation dataset, immune cell proportions as determined by CIBERSORTx [25]. The proteomic discovery and validation datasets were generated using different platforms. Therefore, each dataset was scaled so that all abundance values were between 0 and 1, and then the two datasets were quantile normalised together. Proteins with no missing values that were also found in the proteomic validation dataset were used to train the proteomic classifier.

The DB classifiers were trained to identify DB patients from DV and HC patients, whereas the DV classifiers were trained to identify DV patients from DB and HC patients. Lasso regularised regression [29] was used to identify the discriminatory features and their weights for each classifier. For each sample, a disease risk score (DRS) was calculated using the abundance of the features selected by Lasso, as described by Kaforou et al. in [16]. The DRS was calculated by totalling the abundance of features with positive Lasso weights and subtracting from this total the abundance of features with negative Lasso weights. Features were only included in the DRS if their Lasso weight direction and log-fold change direction were concordant. DRS were scaled between 0 and 1.

The classifiers were tested on the DB and DV patients of their respective validation dataset. The cut-off threshold above which a sample was classified as DB or DV was calculated using the *coords* function in the R package pROC [69] using a sensitivity cut-off of 90%. The classifiers were then tested on the KD patients from the discovery and validation datasets and the thresholds identified by pROC were used to determine if the KD patients were classified as DB or DV. If patients were classified as neither bacterial nor

viral according to their DRS, differential abundance analysis followed by pathway analysis (as described in Section Pathway Analysis) was done to identify the pathways enriched in these patients compared to healthy controls.

#### 4.3.4. Exploration of Kawasaki Disease Samples Alone

In order to identify the natural clusters formed by KD patients in the absence of bacterial or viral patients, *K*-Means clustering was done separately on the KD patients. The process followed was the same as outlined in Section Clustering Analysis. The association between cluster membership and clinical variables was tested. The clinical variables and the statistical tests used were the same as outlined in Section Clustering Analysis.

## 5. Conclusions

Taken together, the results from differential abundance analysis, pathway analysis, clustering and classification suggest that the host transcriptomic and proteomic responses during KD are highly heterogeneous. Different clusters of host responses during KD were identified, some of which resemble elements of host responses to bacterial, viral and fungal infections. These differences in the host responses could imply that KD is triggered either by several different pathogens, or by a single pathogen that has different manifestations according to the underlying genetic and environmental situation of the host. Whilst there are similarities between the host response during KD and the host response to bacterial infections and viral infections, there are also many differences in the responses, suggesting that KD may be triggered by a novel process not typical of either common bacterial or viral infections. This was demonstrated by the majority of the KD transcriptomic samples falling into a non-bacterial, non-viral group following classification, raising the possibility that the minority of KD transcriptomic samples with bacterial or viral profiles were possibly suffering from intercurrent infection in addition to a separate KD trigger. Our data further suggests that research into the aetiologies of KD should be focused on cohorts of KD patients who share similar clinical characteristics in order to identify shared molecular responses.

**Supplementary Materials:** Supplementary materials can be found at <https://www.mdpi.com/article/10.3390/ijms22115655/s1>.

**Author Contributions:** Conceptualisation, M.K. and M.L.; methodology, S.H., H.J., A.M., S.M., C.C.-P., A.T. (Adam Thorne), H.H., R.F., and M.K.; formal analysis, H.J., A.M. and M.K.; investigation, S.H., M.I.d.J., S.M., C.C.-P., C.S., A.T., R.F., and V.W.; resources, J.C.B., R.G., S.H., M.I.d.J., J.K., S.M., C.S., A.T. (Adriana Tremoulet), R.F., A.T. (Adam Thorne) and V.W.; data curation, J.C.B., S.H., J.H., H.J., A.M., S.M., C.S. and V.W.; writing—original draft preparation, H.J.; writing—review and editing, all authors; validation, H.J.; visualisation, H.J.; supervision, J.H., M.L. and M.K.; project administration, H.J., C.S., V.W. and M.K.; funding acquisition, H.J., S.M., S.H., R.G., A.T. (Adriana Tremoulet), M.I.d.J., T.K., V.W., J.C.B., C.C.-P., J.H., M.L. and M.K. All authors have read and agreed to the published version of the manuscript.

**Funding:** H.J. receives support from the Wellcome Trust (4-year PhD programme, grant number 215214/Z/19/Z). M.K. receives support from the Wellcome Trust (Sir Henry Wellcome Fellowship grant number 206508/Z/17/Z). R.G., J.H., S.H., M.I.d.J., M.K., T.K., M.L., S.M., C.C.P. and V.W. acknowledge funding for the EUCLIDS and PERFORM studies, funded by the European Union, grant numbers 668303 and 279185. J.C.B. and A.T. acknowledge the Marilyn and Gordon Macklin Foundation and the National Institutes of Health, Heart, Lung Blood Institute (NHLBI) 1R01 HL140898. J.H., M.L. and M.K. acknowledge support from the NIHR Imperial College BRC.

**Institutional Review Board Statement:** The IRIS, EUCLIDS, PERFORM and UK Kawasaki studies were conducted according to the guidelines of the Declaration of Helsinki, and approved by St Mary's Research Ethics Committee (REC 09/H0712/50, 17/07/2009; REC 11/LO/1982, 09/01/2012; REC 16/LO/1684, 09/11/2016; REC 13/LO/0026, 25/02/2013, respectively). The GENDRES study was conducted according to the guidelines of the Declaration of Helsinki and approved by the Ethical Committee of Clinical Investigation of Galicia (CEIC ref 2010/015, 25/03/2010). The UCSD KD

study was conducted according to the guidelines of the Declaration of Helsinki, and approved by the University of California San Diego (Human Research Protection Program #140220).

**Informed Consent Statement:** Informed consent was obtained from all subjects involved in the study.

**Data Availability Statement:** Transcriptomic datasets used in this study can be found on the Gene Expression Omnibus under accession codes: GSE73461, GSE73462 and GSE73463. Further data may be available on request to the authors.

**Acknowledgments:** Many thanks to the patients and their families who took part in the studies that the data presented here originated from. This study was supported by all members of Personalised Risk assessment in febrile illness to optimise Real-life Management (PERFORM).

**Conflicts of Interest:** The authors declare no conflict of interest. The funders had no role in the design of the study; in the collection, analyses, or interpretation of data; in the writing of the manuscript, or in the decision to publish the results.

## Abbreviations

CAA	Coronary artery aneurysm
DB	Definite bacterial
DRS	Disease risk score
DV	Definite viral
FDR	False discovery rate
HC	Healthy control
KD	Kawasaki disease
LFC	Log-fold change
SDA	Significantly differentially abundant

## References

1. Kawasaki, T.; Kosaki, F.; Okawa, S.; Shigematsu, I.; Yanagawa, H. A New Infantile Acute Febrile Mucocutaneous Lymph Node Syndrome (MLNS) Prevailing in Japan. *Pediatrics* **1974**, *54*.
2. Ramphul, K.; Mejjias, S.G. Kawasaki disease: A comprehensive review. *Arch. Med. Sci. Atheroscler. Dis.* **2018**, *3*, 41–45. [[CrossRef](#)] [[PubMed](#)]
3. Ogata, S.; Shimizu, C.; Franco, A.; Touma, R.; Kanegaye, J.T.; Choudhury, B.P.; Naidu, N.N.; Kanda, Y.; Hoang, L.T.; Hibberd, M.L.; et al. Treatment Response in Kawasaki Disease Is Associated with Sialylation Levels of Endogenous but Not Therapeutic Intravenous Immunoglobulin G. *PLoS ONE* **2013**, *8*, e81448. [[CrossRef](#)] [[PubMed](#)]
4. Skochko, S.M.; Jain, S.; Sun, X.; Sivilyay, N.; Kanegaye, J.T.; Pancheri, J.; Shimizu, C.; Sheets, R.; Tremoulet, A.H.; Burns, J.C. Kawasaki Disease Outcomes and Response to Therapy in a Multiethnic Community: A 10-Year Experience. *J. Pediatr.* **2018**, *203*, 408–415. [[CrossRef](#)] [[PubMed](#)]
5. Brogan, P.; Burns, J.C.; Cornish, J.; Diwakar, V.; Eleftheriou, D.; Gordon, J.B.; Gray, H.H.; Johnson, T.W.; Levin, M.; Malik, I.; et al. Lifetime cardiovascular management of patients with previous Kawasaki disease. *Heart* **2020**, *106*, 411–420. [[CrossRef](#)]
6. Singh, S.; Vignesh, P.; Burgner, D. The epidemiology of Kawasaki disease: A global update. *Arch. Dis. Child.* **2015**, *100*, 1084–1088. [[CrossRef](#)]
7. Nagata, S. Causes of Kawasaki Disease—From Past to Present. *Front. Pediatr.* **2019**, *7*, 18. [[CrossRef](#)]
8. Dietz, S.M.; van Stijn, D.; Burgner, D.; Levin, M.; Kuipers, I.M.; Hutten, B.A.; Kuijpers, T.W. Dissecting Kawasaki disease: A state-of-the-art review. *Eur. J. Pediatr.* **2017**, *176*, 995–1009. [[CrossRef](#)]
9. Nakamura, A.; Ikeda, K.; Hamaoka, K. Aetiological significance of infectious stimuli in Kawasaki disease. *Front. Pediatr.* **2019**, *7*, 244. [[CrossRef](#)]
10. Rodó, X.; Ballester, J.; Cayan, D.; Melish, M.E.; Nakamura, Y.; Uehara, R.; Burns, J.C. Association of Kawasaki disease with tropospheric wind patterns. *Sci. Rep.* **2011**, *1*, 150. [[CrossRef](#)]
11. Rypdal, M.; Rypdal, V.; Burney, J.A.; Cayan, D.; Bainto, E.; Skochko, S.; Tremoulet, A.H.; Creamean, J.; Shimizu, C.; Kim, J.; et al. Clustering and climate associations of Kawasaki Disease in San Diego County suggest environmental triggers. *Sci. Rep.* **2018**, *8*, 1–9. [[CrossRef](#)]
12. Levin, M. Childhood Multisystem Inflammatory Syndrome—A New Challenge in the Pandemic. *N. Engl. J. Med.* **2020**, *383*, 393–395. [[CrossRef](#)]
13. Whittaker, E.; Bamford, A.; Kenny, J.; Kaforou, M.; Jones, C.E.; Shah, P.; Ramnarayan, P.; Fraisse, A.; Miller, O.; Davies, P.; et al. Clinical Characteristics of 58 Children with a Pediatric Inflammatory Multisystem Syndrome Temporally Associated with SARS-CoV-2. *JAMA J. Am. Med. Assoc.* **2020**, *324*, 259–269. [[CrossRef](#)]

14. Dufort, E.M.; Koumans, E.H.; Chow, E.J.; Rosenthal, E.M.; Muse, A.; Rowlands, J.; Barranco, M.A.; Maxted, A.M.; Rosenberg, E.S.; Easton, D.; et al. Multisystem Inflammatory Syndrome in Children in New York State. *N. Engl. J. Med.* **2020**, *383*, 347–358. [[CrossRef](#)]
15. McCrindle, B.W.; Manlhiot, C. SARS-CoV-2-Related Inflammatory Multisystem Syndrome in Children: Different or Shared Etiology and Pathophysiology as Kawasaki Disease? *JAMA J. Am. Med. Assoc.* **2020**, *324*, 246–248. [[CrossRef](#)]
16. Kaforou, M.; Wright, V.J.; Oni, T.; French, N.; Anderson, S.T.; Bangani, N.; Banwell, C.M.; Brent, A.J.; Crampin, A.C.; Dockrell, H.M.; et al. Detection of Tuberculosis in HIV-Infected and -Uninfected African Adults Using Whole Blood RNA Expression Signatures: A Case-Control Study. *PLoS Med.* **2013**, *10*, e1001538. [[CrossRef](#)]
17. Kaforou, M.; Herberg, J.A.; Wright, V.J.; Coin, L.J.M.; Levin, M. Diagnosis of Bacterial Infection Using a 2-Transcript Host RNA Signature in Febrile Infants 60 Days or Younger. *JAMA* **2017**, *317*, 1577–1578. [[CrossRef](#)]
18. Wright, V.J.; Herberg, J.A.; Kaforou, M.; Shimizu, C.; Eleftherohorinou, H.; Shailes, H.; Barendregt, A.M.; Menikou, S.; Gormley, S.; Berk, M.; et al. Diagnosis of Kawasaki Disease Using a Minimal Whole-Blood Gene Expression Signature. *JAMA Pediatr.* **2018**, *172*, e182293. [[CrossRef](#)]
19. Hoang, L.T.; Shimizu, C.; Ling, L.; Naim, A.N.M.; Khor, C.C.; Tremoulet, A.H.; Wright, V.; Levin, M.; Hibberd, M.L.; Burns, J.C. Global gene expression profiling identifies new therapeutic targets in acute Kawasaki disease. *Genome Med.* **2014**. [[CrossRef](#)]
20. Gold, L.; Ayers, D.; Bertino, J.; Bock, C.; Bock, A.; Brody, E.N.; Carter, J.; Dalby, A.B.; Eaton, B.E.; Fitzwater, T.; et al. Aptamer-based multiplexed proteomic technology for biomarker discovery. *PLoS ONE* **2010**. [[CrossRef](#)]
21. McCrindle, B.W.; Rowley, A.H.; Newburger, J.W.; Burns, J.C.; Bolger, A.F.; Gewitz, M.; Baker, A.L.; Jackson, M.A.; Takahashi, M.; Shah, P.B.; et al. Diagnosis, treatment, and long-term management of Kawasaki disease: A scientific statement for health professionals from the American Heart Association. *Circulation* **2017**, *135*, e927–e999. [[CrossRef](#)]
22. Herberg, J.A.; Kaforou, M.; Gormley, S.; Sumner, E.R.; Patel, S.; Jones, K.D.J.; Paulus, S.; Fink, C.; Martinon-Torres, F.; Montana, G.; et al. Transcriptomic profiling in childhood H1N1/09 influenza reveals reduced expression of protein synthesis genes. *J. Infect. Dis.* **2013**, *208*, 1664–1668. [[CrossRef](#)]
23. Ritchie, M.E.; Phipson, B.; Wu, D.; Hu, Y.; Law, C.W.; Shi, W.; Smyth, G.K. limma powers differential expression analyses for RNA-sequencing and microarray studies. *Nucleic Acids Res.* **2015**, *43*, e47. [[CrossRef](#)]
24. Raudvere, U.; Kolberg, L.; Kuzmin, I.; Arak, T.; Adler, P.; Peterson, H.; Vilo, J. g:Profiler: a web server for functional enrichment analysis and conversions of gene lists (2019 update). *Nucleic Acids Res.* **2019**, *47*, W191–W198. [[CrossRef](#)]
25. Newman, A.M.; Steen, C.B.; Liu, C.L.; Gentles, A.J.; Chaudhuri, A.A.; Scherer, F.; Khodadoust, M.S.; Esfahani, M.S.; Luca, B.A.; Steiner, D.; et al. Determining cell type abundance and expression from bulk tissues with digital cytometry. *Nat. Biotechnol.* **2019**, *37*, 773–782. [[CrossRef](#)]
26. Charrad, M.; Ghazzali, N.; Boiteau, V.; Niknafs, A. Nbclust: An R package for determining the relevant number of clusters in a data set. *J. Stat. Softw.* **2014**. [[CrossRef](#)]
27. Kimura, Y.; Yanagimachi, M.; Ino, Y.; Aketagawa, M.; Matsuo, M.; Okayama, A.; Shimizu, H.; Oba, K.; Morioka, I.; Imagawa, T.; et al. Identification of candidate diagnostic serum biomarkers for Kawasaki disease using proteomic analysis. *Sci. Rep.* **2017**. [[CrossRef](#)] [[PubMed](#)]
28. Whitin, J.C.; Yu, T.T.S.; Ling, X.B.; Kanegaye, J.T.; Burns, J.C.; Cohen, H.J. A novel truncated form of serum amyloid a in kawasaki disease. *PLoS ONE* **2016**, *11*. [[CrossRef](#)] [[PubMed](#)]
29. Tibshirani, R. Regression Shrinkage and Selection via the Lasso. *J. R. Stat. Soc. Ser. B* **1996**, *58*, 267–288. [[CrossRef](#)]
30. Danese, S. Nonimmune cells in inflammatory bowel disease: From victim to villain. *Trends Immunol.* **2008**, *29*, 555–564. [[CrossRef](#)]
31. Tanaka, K.; Yoshioka, T.; Bieberich, C.; Jay, G. Role of the Major Histocompatibility Complex Class I Antigens in Tumor Growth and Metastasis. *Annu. Rev. Immunol.* **1988**, *6*, 359–380. [[CrossRef](#)]
32. Tremoulet, A.H.; Jain, S.; Chandrasekar, D.; Sun, X.; Sato, Y.; Burns, J.C. Evolution of laboratory values in patients with Kawasaki disease. *Pediatr. Infect. Dis. J.* **2011**, *30*, 1022–1026. [[CrossRef](#)]
33. Biezeveld, M.H.; van Mierlo, G.; Lutter, R.; Kuipers, I.M.; Dekker, T.; Hack, C.E.; Newburger, J.W.; Kuijpers, T.W. Sustained activation of neutrophils in the course of Kawasaki disease: An association with matrix metalloproteinases. *Clin. Exp. Immunol.* **2005**, *141*, 183–188. [[CrossRef](#)]
34. Asano, T.; Ogawa, S. Expression of IL-8 in Kawasaki disease. *Clin. Exp. Immunol.* **2000**, *122*, 514–519. [[CrossRef](#)]
35. Zandstra, J.; van de Geer, A.; Tanck, M.W.T.; van Stijn-Bringas Dimitriades, D.; Aarts, C.E.M.; Dietz, S.M.; van Bruggen, R.; Schweintzger, N.A.; Zenz, W.; Emonts, M.; et al. Biomarkers for the Discrimination of Acute Kawasaki Disease From Infections in Childhood. *Front. Pediatr.* **2020**, *8*, 355. [[CrossRef](#)]
36. Manlhiot, C.; Mueller, B.; O’Shea, S.; Majeed, H.; Bernknopf, B.; Labelle, M.; Westcott, K.V.; Bai, H.; Chahal, N.; Birken, C.S.; et al. Environmental epidemiology of Kawasaki disease: Linking disease etiology, pathogenesis and global distribution. *PLoS ONE* **2018**, *13*. [[CrossRef](#)]
37. Rodó, X.; Curcoll, R.; Robinson, M.; Ballester, J.; Burns, J.C.; Cayan, D.R.; Lipkin, W.I.; Williams, B.L.; Couto-Rodriguez, M.; Nakamura, Y.; et al. Tropospheric winds from northeastern China carry the etiologic agent of Kawasaki disease from its source to Japan. *Proc. Natl. Acad. Sci. USA* **2014**, *111*, 7952–7957. [[CrossRef](#)]
38. Lu, X.; Zhang, L.; Du, H.; Zhang, J.; Li, Y.Y.; Qu, J.; Zhang, W.; Wang, Y.; Bao, S.; Li, Y.; et al. SARS-CoV-2 Infection in Children. *N. Engl. J. Med.* **2020**, *382*, 1663–1665. [[CrossRef](#)]

39. Göttinger, F.; Santiago-García, B.; Noguera-Julián, A.; Lanaspá, M.; Lancella, L.; Calò Carducci, F.I.; Gabrovská, N.; Velizarova, S.; Prunk, P.; Osterman, V.; et al. COVID-19 in children and adolescents in Europe: A multinational, multicentre cohort study. *Lancet Child Adolesc. Heal.* **2020**, *0*. [CrossRef]
40. Davies, P.; Evans, C.; Kanthimathinathan, H.K.; Lillie, J.; Brierley, J.; Waters, G.; Johnson, M.; Griffiths, B.; du Pré, P.; Mohammad, Z.; et al. Intensive care admissions of children with paediatric inflammatory multisystem syndrome temporally associated with SARS-CoV-2 (PIMS-TS) in the UK: a multicentre observational study. *Lancet Child Adolesc. Health* **2020**. [CrossRef]
41. Burgner, D.; Davila, S.; Breunis, W.B.; Ng, S.B.; Li, Y.; Bonnard, C.; Ling, L.; Wright, V.J.; Thalamuthu, A.; Odam, M.; et al. A genome-wide association study identifies novel and functionally related susceptibility loci for Kawasaki disease. *PLoS Genet.* **2009**. [CrossRef] [PubMed]
42. Curtis, N.; Zheng, R.; Lamb, J.R.; Levin, M. Evidence for a superantigen mediated process in Kawasaki disease. *Arch. Dis. Child.* **1995**, *72*, 308–311. [CrossRef] [PubMed]
43. Han, S.B.; Lee, S.Y. Antibiotic use in children with Kawasaki disease. *World J. Pediatr.* **2018**, *14*, 621–622. [CrossRef] [PubMed]
44. Diz, A.P.; Martínez-Fernández, M.; Rolán-Alvarez, E. Proteomics in evolutionary ecology: Linking the genotype with the phenotype. *Mol. Ecol.* **2012**, *21*, 1060–1080. [CrossRef]
45. Benseler, S.M.; McCrindle, B.W.; Silverman, E.D.; Tyrrell, P.N.; Wong, J.; Yeung, R.S.M. Infections and Kawasaki disease: Implications for coronary artery outcome. *Pediatrics* **2005**, *116*, e760–e766. [CrossRef]
46. Martínón-Torres, F.; Salas, A.; Rivero-Calle, I.; Cebejón-López, M.; Pardo-Seco, J.; Herberg, J.A.; Boeddha, N.P.; Klobassa, D.S.; Secka, F.; Paulus, S.; et al. Life-threatening infections in children in Europe (the EUCLIDS Project): a prospective cohort study. *Lancet Child Adolesc. Health* **2018**, *2*, 404–414. [CrossRef]
47. R Foundation for Statistical Computing R: A language and environment for statistical computing. *R A Lang. Environ. Stat. Comput.* **3.3.1** **2016**.
48. Du, P.; Kibbe, W.A.; Lin, S.M. Lumi: A pipeline for processing Illumina microarray. *Bioinformatics* **2008**. [CrossRef]
49. Leek, J.T.; Johnson, W.E.; Parker, H.S.; Jaffe, A.E.; Storey, J.D. The SVA package for removing batch effects and other unwanted variation in high-throughput experiments. *Bioinformatics* **2012**, *28*, 882–883. [CrossRef]
50. Cox, J.; Mann, M. MaxQuant enables high peptide identification rates, individualized p.p.b.-range mass accuracies and proteome-wide protein quantification. *Nat. Biotechnol.* **2008**, *26*, 1367–1372. [CrossRef]
51. Cox, J.; Hein, M.Y.; Lubner, C.A.; Paron, I.; Nagaraj, N.; Mann, M. Accurate proteome-wide label-free quantification by delayed normalization and maximal peptide ratio extraction, termed MaxLFQ. *Mol. Cell. Proteomics* **2014**, *13*, 2513–2526. [CrossRef]
52. Sweeney, T.E.; COCONUT: COmbat CO-Normalization Using conTrols (COCONUT). R Package Version 1.0.2. 2017. Available online: <https://rdrr.io/cran/COCONUT/> (accessed on 25 May 2021).
53. Benjamini, Y.; Hochberg, Y. Controlling the False Discovery Rate: A Practical and Powerful Approach to Multiple Testing. *J. R. Stat. Soc. Ser. B* **1995**, *57*, 289–300. [CrossRef]
54. Supek, F.; Bošnjak, M.; Škunca, N.; Šmuc, T. REVIGO Summarizes and Visualizes Long Lists of Gene Ontology Terms. *PLoS ONE* **2011**, *6*, e21800. [CrossRef]
55. Hartigan, J.A.; Wong, M.A. Algorithm AS 136: A K-Means Clustering Algorithm. *Appl. Stat.* **1979**, *28*, 100. [CrossRef]
56. Krzanowski, W.J.; Lai, Y.T. A Criterion for Determining the Number of Groups in a Data Set Using Sum-of-Squares Clustering. *Biometrics* **1988**, *44*, 23. [CrossRef]
57. Caliński, T.; Harabasz, J. A Dendrite Method For Cluster Analysis. *Commun. Stat.* **1974**, *3*, 1–27. [CrossRef]
58. Gordon, A.D.; Hartigan, J.A. Clustering Algorithms. *J. Am. Stat. Assoc.* **1976**. [CrossRef]
59. McClain, J.O.; Rao, V.R. CLUSTISZ: A Program to Test for the Quality of Clustering of a Set of Objects. *J. Mark. Res.* **1975**, *12*, 456–460.
60. Dunn, J.C. Well-separated clusters and optimal fuzzy partitions. *J. Cybern.* **1974**, *4*, 95–104. [CrossRef]
61. Halkidi, M.; Vazirgiannis, M.; Balisalakos, V. Quality scheme assessment in the clustering process. In *Principles of Data Mining and Knowledge Discovery*; Lecture Notes in Computer Science (including subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics); Springer: Berlin/Heidelberg, Germany, 2000; Volume 1910, pp. 265–276.
62. Halkidi, M.; Vazirgiannis, M. Clustering validity assessment: Finding the optimal partitioning of a data set. In *Proceedings of the IEEE International Conference on Data Mining, ICDM, San Jose, CA, USA, 29 November–2 December 2001*; pp. 187–194. [CrossRef]
63. Hubert, L.J.; Levin, J.R. A general statistical framework for assessing categorical clustering in free recall. *Psychol. Bull.* **1976**, *83*, 1072–1080. [CrossRef]
64. Rousseeuw, P.J. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *J. Comput. Appl. Math.* **1987**, *20*, 53–65. [CrossRef]
65. Ball, G.H.; Hall, D.J. *Isodata, A Novel Method of Data Analysis and Pattern Classification*; Stanford Research Institute: Menlo Park, CA, USA, 1965.
66. Milligan, G.W. An examination of the effect of six types of error perturbation on fifteen clustering algorithms. *Psychometrika* **1980**, *45*, 325–342. [CrossRef]
67. Milligan, G.W. A monte carlo study of thirty internal criterion measures for cluster analysis. *Psychometrika* **1981**, *46*, 187–199. [CrossRef]



- 
68. Ratkowsky, D.; Lance, G. A Criterion for Determining the Number of Groups in a Classification. *Aust. Comput. J.* **1978**, *10*, 115–117.
  69. Robin, X.; Turck, N.; Hainard, A.; Tiberti, N.; Lisacek, F.; Sanchez, J.C.; Müller, M. pROC: An open-source package for R and S+ to analyze and compare ROC curves. *BMC Bioinform.* **2011**, *12*, 77. [[CrossRef](#)]