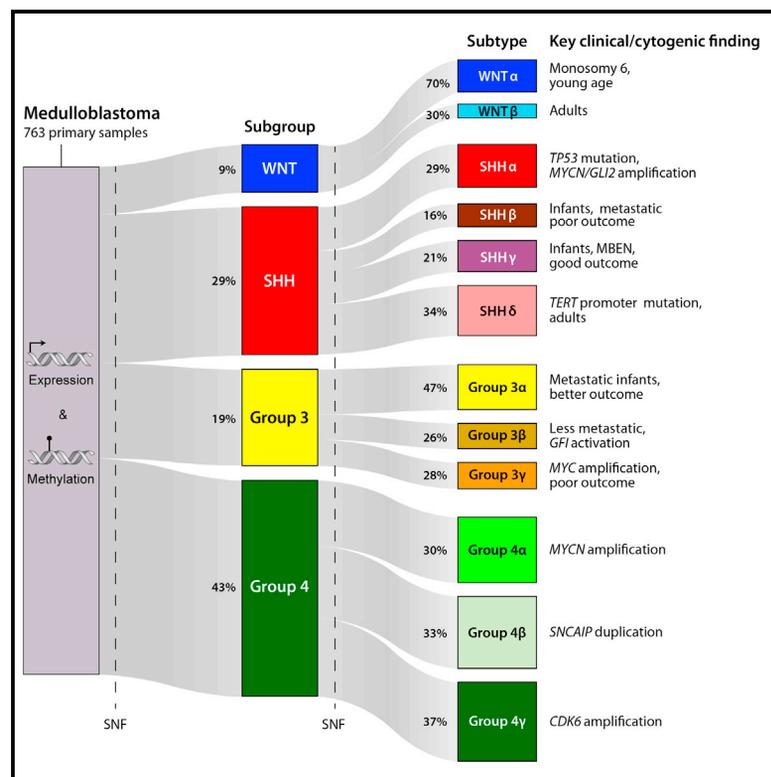


Intertumoral Heterogeneity within Medulloblastoma Subgroups

Graphical Abstract



Authors

Florence M.G. Cavalli, Marc Remke, Ladislav Rampasek, ..., Anna Goldenberg, Vijay Ramaswamy, Michael D. Taylor

Correspondence

anna.goldenberg@utoronto.ca (A.G.), vijay.ramaswamy@sickkids.ca (V.R.), mdtaylor@sickkids.ca (M.D.T.)

In Brief

Cavalli et al. analyze 763 primary medulloblastoma samples using the similarity network fusion approach. They identify subtypes that have distinct somatic copy-number aberrations, activated pathways, and clinical outcomes within each of the four known subgroups and further delineate group 3 from group 4 MB.

Highlights

- Medulloblastoma comprises 12 subtypes; 2 WNT, 4 SHH, 3 group 3, and 3 group 4 groups
- Heterogeneity within subgroups accounts for previously unexplained variation
- Groups 3 and 4 medulloblastoma are molecularly distinct entities
- Clinically and biologically relevant subtypes exist for each subgroup



Intertumoral Heterogeneity within Medulloblastoma Subgroups

Florence M.G. Cavalli,^{1,2,80} Marc Remke,^{3,4,71,80} Ladislav Rampasek,^{5,6} John Peacock,^{1,2,4} David J.H. Shih,^{1,2,4} Betty Luu,^{1,2} Livia Garzia,^{1,2} Jonathon Torchia,^{1,4} Carolina Nor,^{1,2} A. Sorana Morrissy,^{1,2} Sameer Agnihotri,⁷ Yuan Yao Thompson,^{1,2,4} Claudia M. Kuzan-Fischer,^{1,2} Hamza Farooq,^{1,2,4} Keren Isaev,^{8,9} Craig Daniels,^{1,2} Byung-Kyu Cho,¹⁰ Seung-Ki Kim,¹⁰ Kyu-Chang Wang,¹⁰ Ji Yeoun Lee,¹⁰ Wieslawa A. Grajkowska,¹¹

(Author list continued on next page)

¹The Arthur and Sonia Labatt Brain Tumour Research Centre

²Developmental & Stem Cell Biology Program

The Hospital for Sick Children, Toronto, ON M5G 1X8, Canada

³Department of Pediatric Oncology, Hematology, and Clinical Immunology, Medical Faculty, University Hospital Düsseldorf, Düsseldorf 40225, Germany

⁴Department of Laboratory Medicine and Pathobiology, University of Toronto, Toronto, ON M5S 1A1, Canada

⁵Department of Computer Science, University of Toronto, Toronto, ON M5S 2E4, Canada

⁶Program in Genetics and Genome Biology, The Hospital for Sick Children, Toronto, ON M5G 1X8, Canada

⁷UPCI Brain Tumor Program, University of Pittsburgh, Children's Hospital of Pittsburgh, Pittsburgh, PA 15224, USA

⁸Informatics Program, Ontario Institute for Cancer Research, Toronto, ON M5G 0A3, Canada

⁹Department of Medical Biophysics, University of Toronto, Toronto, ON M5G 1L7, Canada

¹⁰Department of Neurosurgery, Division of Pediatric Neurosurgery, Seoul National University Children's Hospital, Seoul 30322, South Korea

¹¹Department of Pathology

(Affiliations continued on next page)

SUMMARY

While molecular subgrouping has revolutionized medulloblastoma classification, the extent of heterogeneity within subgroups is unknown. Similarity network fusion (SNF) applied to genome-wide DNA methylation and gene expression data across 763 primary samples identifies very homogeneous clusters of patients, supporting the presence of medulloblastoma subtypes. After integration of somatic copy-number alterations, and clinical features specific to each cluster, we identify 12 different subtypes of medulloblastoma. Integrative analysis using SNF further delineates group 3 from group 4 medulloblastoma, which is not as readily apparent through analyses of individual data types. Two clear subtypes of infants with Sonic Hedgehog medulloblastoma with disparate outcomes and biology are identified. Medulloblastoma subtypes identified through integrative clustering have important implications for stratification of future clinical trials.

INTRODUCTION

Genomics has substantially advanced our understanding of medulloblastoma (Northcott et al., 2012a; Ramaswamy et al.,

2011). While historically considered one entity, it is now clearly accepted that medulloblastoma comprises at least four distinct entities: WNT, SHH, group 3, and group 4; as reflected in the current revision of the WHO classification (Louis et al., 2016;

Significance

While medulloblastoma is widely recognized as comprising four distinct subgroups, the degree of heterogeneity within the four subgroups, and the extent of overlap between the four subgroups is unknown. Applying similarity network fusion to integrate gene expression and DNA methylation profiling, we demonstrate that the degree of overlap between groups 3 and 4 is minimal after accounting for both expression and DNA methylation data. We identify medulloblastoma subtypes within each of the subgroups that have distinct somatic copy-number aberrations, differentially activated pathways, and disparate clinical outcomes. Integrated analysis has refined the boundaries between the four medulloblastoma subgroups, and identified clinically and biologically relevant subtypes, which will inform and improve preclinical modeling, as well as refine our current clinical classification.

Marta Perek-Polnik,¹² Alexandre Vasiljevic,^{13,72} Cecile Faure-Conter,¹⁴ Anne Jouvett,¹⁵ Caterina Giannini,¹⁶ Amulya A. Nageswara Rao,¹⁷ Kay Ka Wai Li,¹⁸ Ho-Keung Ng,¹⁸ Charles G. Eberhart,¹⁹ Ian F. Pollack,²⁰ Ronald L. Hamilton,²¹ G. Yancey Gillespie,²² James M. Olson,^{23,24} Sarah Leary,²⁴ William A. Weiss,²⁵ Boleslaw Lach,^{26,73} Lola B. Chambless,²⁷ Reid C. Thompson,²⁷ Michael K. Cooper,²⁸ Rajeev Vibhakar,²⁹ Peter Hauser,³⁰ Marie-Lise C. van Veelen,³¹ Johan M. Kros,³² Pim J. French,³³ Young Shin Ra,³⁴ Toshihiro Kumabe,³⁵ Enrique López-Aguilar,³⁶ Karel Zitterbart,³⁷ Jaroslav Sterba,³⁷ Gaetano Finocchiaro,³⁸ Maura Massimino,³⁹ Erwin G. Van Meir,⁴⁰ Satoru Osuka,⁴⁰ Tomoko Shofuda,⁴¹ Almos Klekner,⁴² Massimo Zollo,⁴³ Jeffrey R. Leonard,⁴⁴ Joshua B. Rubin,⁴⁵ Nada Jabado,⁴⁶ Steffen Albrecht,^{47,74} Jaume Mora,⁴⁸ Timothy E. Van Meter,⁴⁹ Shin Jung,⁵⁰

(Author list continued on next page)

¹²Department of Oncology

The Children's Memorial Health Institute, University of Warsaw, Warsaw 04-730, Poland

¹³Centre de Pathologie et Neuropathologie Est, Centre de Biologie et Pathologie Est, Groupement Hospitalier Est, Hospices Civils de Lyon, Bron 69677, France

¹⁴Institute of Pediatric Hematology and Oncology, Lyon 69008, France

¹⁵Centre de Pathologie EST, Groupement Hospitalier EST, Université de Lyon, Bron 69677, France

¹⁶Department of Laboratory Medicine and Pathology

¹⁷Division of Pediatric Hematology/Oncology

Mayo Clinic, Rochester, MN 55905, USA

¹⁸Department of Anatomical and Cellular Pathology, The Chinese University of Hong Kong, Shatin, New Territories, Hong Kong, China

¹⁹Departments of Pathology, Ophthalmology and Oncology, John Hopkins University School of Medicine, Baltimore, MD 21287, USA

²⁰Department of Neurological Surgery, University of Pittsburgh School of Medicine, Pittsburgh, PA 15213, USA

²¹Department of Pathology, University of Pittsburgh School of Medicine, Pittsburgh, PA 15261, USA

²²Department of Surgery, Division of Neurosurgery, University of Alabama at Birmingham, Birmingham, AL 35233, USA

²³Clinical Research Division, Fred Hutchinson Cancer Research Center, Seattle, WA 98109-1024, USA

²⁴Division of Pediatric Hematology/Oncology, University of Washington School of Medicine, Seattle Children's Hospital, Seattle, WA 98145-5005, USA

²⁵Departments of Pediatrics, Neurological Surgery and Neurology, University of California San Francisco, San Francisco, CA 94143-0112, USA

²⁶Division of Anatomical Pathology, Department of Pathology and Molecular Medicine, McMaster University, Hamilton, ON L8S 4K1, Canada

²⁷Department of Neurological Surgery

²⁸Department of Neurology

Vanderbilt Medical Center, Nashville, TN 37232, USA

²⁹Department of Pediatrics, University of Colorado Denver, Aurora, CO 80045, USA

³⁰2nd Department of Pediatrics, Semmelweis University, Budapest 1094, Hungary

³¹Department of Neurosurgery, Erasmus University Medical Center, Rotterdam 3015 CE, the Netherlands

³²Department of Pathology, Erasmus University Medical Center, Rotterdam 3015 CN, the Netherlands

³³Department of Neurology, Erasmus University Medical Center, Rotterdam 3015 CE, the Netherlands

³⁴Department of Neurosurgery, University of Ulsan, Asan Medical Center, Seoul 05505, South Korea

³⁵Department of Neurosurgery, Kitasato University School of Medicine, Sagamihara, Kanagawa 252-0374, Japan

³⁶Division of Pediatric Hematology/Oncology, Hospital Pediatría Centro Médico Nacional Century XXI, Mexico City 06720, Mexico

³⁷Department of Pediatric Oncology, School of Medicine, Masaryk University, Brno 625 00, Czech Republic

³⁸Department of Neuro-Oncology, Istituto Neurologico Besta

³⁹Fondazione IRCCS Istituto Nazionale Tumori

Milan 20133, Italy

⁴⁰Department of Hematology & Medical Oncology, School of Medicine and Winship Cancer Institute, Emory University, Atlanta, GA 30322, USA

⁴¹Division of Stem Cell Research, Institute for Clinical Research, Osaka National Hospital, Osaka 540-0006, Japan

⁴²Department of Neurosurgery, University of Debrecen, Medical and Health Science Centre, Debrecen 4032, Hungary

⁴³Dipartimento di Biochimica e Biotecnologie Mediche, University of Naples, Naples 80145, Italy

(Affiliations continued on next page)

Ramaswamy et al., 2016a). These four subgroups have distinct transcriptional profiles, copy-number aberrations, somatic mutations, and clinical outcomes (Morrissy et al., 2016; Northcott et al., 2012a; Ramaswamy et al., 2016b; Ramaswamy et al., 2013). Indeed, current clinical trials and risk stratification biomarkers incorporate the four molecular subgroups (Ramaswamy et al., 2016a), as do preclinical modeling and the development of novel therapeutics (Pei et al., 2016). However, the extent to which there are additional layers of heterogeneity within the me-

dulloblastoma subgroups is unknown, and a concerted global effort to analyze a very large cohort of tumors will be needed to resolve the question.

WNT and SHH medulloblastomas are clearly identifiable and separable across the majority of transcriptional and methylation profiling studies, demonstrating minimal overlap with other subgroups (Taylor et al., 2012). Clear heterogeneity exists within the SHH subgroup, which includes infants, children, and adults, although the extent and nature of the substructure is not clearly

Andrew S. Moore,^{51,75} Andrew R. Hallahan,^{51,75} Jennifer A. Chan,⁵² Daniela P.C. Tirapelli,⁵³ Carlos G. Carlotti,⁵³ Maryam Fouladi,⁵⁴ José Pimentel,⁵⁵ Claudia C. Faria,⁵⁶ Ali G. Saad,⁵⁷ Luca Massimi,⁵⁸ Linda M. Liao,⁵⁹ Helen Wheeler,⁶⁰ Hideo Nakamura,⁶¹ Samer K. Elbabaa,⁶² Mario Perezpeña-Diazconti,⁶³ Fernando Chico Ponce de León,⁶⁴ Shenandoah Robinson,⁶⁵ Michal Zapotocky,⁶⁶ Alvaro Lassaletta,⁶⁶ Annie Huang,^{1,66} Cynthia E. Hawkins,^{1,67} Uri Tabori,^{1,66} Eric Bouffet,^{1,66} Ute Bartels,⁶⁶ Peter B. Dirks,^{1,68} James T. Rutka,^{1,4,68} Gary D. Bader,^{69,76,77,78,79} Jüri Reimand,^{8,9} Anna Goldenberg,^{5,6,*} Vijay Ramaswamy,^{1,66,70,**} and Michael D. Taylor^{1,2,4,68,81,***}

⁴⁴Division of Pediatric Neurosurgery, Department of Neurosurgery

⁴⁵Departments of Pediatrics, Anatomy and Neurobiology

Washington University School of Medicine and St. Louis Children's Hospital, St. Louis, MO 63110, USA

⁴⁶Division of Hematology/Oncology, Department of Pediatrics

⁴⁷Department of Pathology

McGill University, Montreal, QC H4A 3J1, Canada

⁴⁸Developmental Tumor Biology Laboratory, Hospital Sant Joan de Déu, Esplugues de Llobregat, Barcelona 08950, Spain

⁴⁹Department of Pediatrics, Virginia Commonwealth University, School of Medicine, Richmond, VA 23298-0646, USA

⁵⁰Department of Neurosurgery, Chonnam National University Research Institute of Medical Sciences, Chonnam National University Hwasun Hospital and Medical School, Hwasun-gun 519-763, Chonnam South Korea

⁵¹Lady Cilento Children's Hospital, The University of Queensland, Brisbane QLD 4102, Australia

⁵²Department of Pathology and Laboratory Medicine, University of Calgary, Calgary, AB T2N 2T9, Canada

⁵³Department of Surgery and Anatomy, Faculty of Medicine of Ribeirão Preto, University of São Paulo, São Paulo 14049-900, Brazil

⁵⁴Division of Hematology/Oncology, University of Cincinnati, Cincinnati Children's Hospital Medical Center, Cincinnati, OH 45229, USA

⁵⁵Division of Pathology

⁵⁶Division of Neurosurgery

Centro Hospitalar Lisboa Norte, Hospital de Santa Maria, Lisbon 1649-035, Portugal

⁵⁷Department of Pathology, University of Arkansas for Medical Sciences, Little Rock, AR 72205, USA

⁵⁸Department of Pediatric Neurosurgery, Catholic University Medical School, Rome 00198, Italy

⁵⁹Department of Neurosurgery, David Geffen School of Medicine at UCLA, Los Angeles, CA 90095, USA

⁶⁰Kolling Institute of Medical Research, The University of Sydney, Sydney, NSW 2065, Australia

⁶¹Department of Neurosurgery, Kumamoto University Graduate School of Medical Science, Kumamoto 860-8555, Japan

⁶²Division of Pediatric Neurosurgery, Department of Neurosurgery, Saint Louis University School of Medicine, St. Louis, MO, USA

⁶³Department of Pathology

⁶⁴Department of Neurosurgery

Hospital Infantil de Mexico Federico Gomez, Mexico City 06720, Mexico

⁶⁵Division of Pediatric Neurosurgery, Rainbow & Babies Children's Hospital, Case Western Reserve, Cleveland, OH 44106, USA

⁶⁶Division of Haematology / Oncology

⁶⁷Division of Pathology

⁶⁸Division of Neurosurgery

The Hospital for Sick Children, Toronto, ON M5G 1X8, Canada

⁶⁹The Donnelly Centre, University of Toronto, Toronto, ON M5S 3E1, Canada

⁷⁰Program in Neuroscience and Mental Health and Division of Neurology, The Hospital for Sick Children, Toronto, ON M5G 1X8, Canada

⁷¹Department of Pediatric Neuro-Oncogenomics, German Cancer Consortium (DKTK) and German Cancer Research Center (DKFZ), Düsseldorf 40225, Germany

⁷²ONCOFLAM - Neuro-Oncologie et Neuro-Inflammation Centre de Recherche en Neurosciences de Lyon, Lyon 69008, France

⁷³Department of Pathology and Laboratory Medicine, Hamilton General Hospital, Hamilton, ON L8L 2X2, Canada

⁷⁴Department of Pathology, Montreal Children's Hospital, Montreal, QC H4A 3J1, Canada

⁷⁵Oncology Service, Children's Health Queensland Hospital and Health Service, South Brisbane, QLD 4029, Australia

⁷⁶Banting and Best Department of Medical Research, University of Toronto, Toronto, ON M5G 1L6, Canada

⁷⁷McLaughlin Centre, University of Toronto, Toronto, ON M5G 0A4, Canada

⁷⁸Department of Molecular Genetics, University of Toronto, Toronto, ON M5S 1A8, Canada

⁷⁹Samuel Lunenfeld Research Institute at Mount Sinai Hospital, University of Toronto, Toronto, ON M5G 1X5, Canada

⁸⁰These authors contributed equally

⁸¹Lead Contact

*Correspondence: anna.goldenberg@utoronto.ca (A.G.), vijay.ramaswamy@sickkids.ca (V.R.), mtdtaylor@sickkids.ca (M.D.T.)

<http://dx.doi.org/10.1016/j.ccell.2017.05.005>

defined (Northcott et al., 2011; Kool et al., 2014; Lafay-Cousin et al., 2016). The transcriptomes of group 3 and group 4 medulloblastoma are more similar to each other, and several cytogenetic features, such as isochromosome 17q (i17q), are found in both groups (Taylor et al., 2012). In response to this, the recent revision of WHO Classification of CNS Tumors has assigned groups 3 and 4 as provisional entities, and a recent consensus on high-risk medulloblastoma left this question unresolved (Louis et al., 2016). Establishing the nature of the bound-

ary between group 3 and group 4 is of clinical importance as outcomes differ, particularly in the setting of upfront metastatic dissemination (Ramaswamy et al., 2016a, 2016b; Thompson et al., 2016).

Genome-wide transcriptional arrays and/or genome-wide methylation arrays are the current gold standard for medulloblastoma subgrouping (Ramaswamy et al., 2016a). These approaches have been used independently with the underlying assumption that they identify similar, perhaps even identical

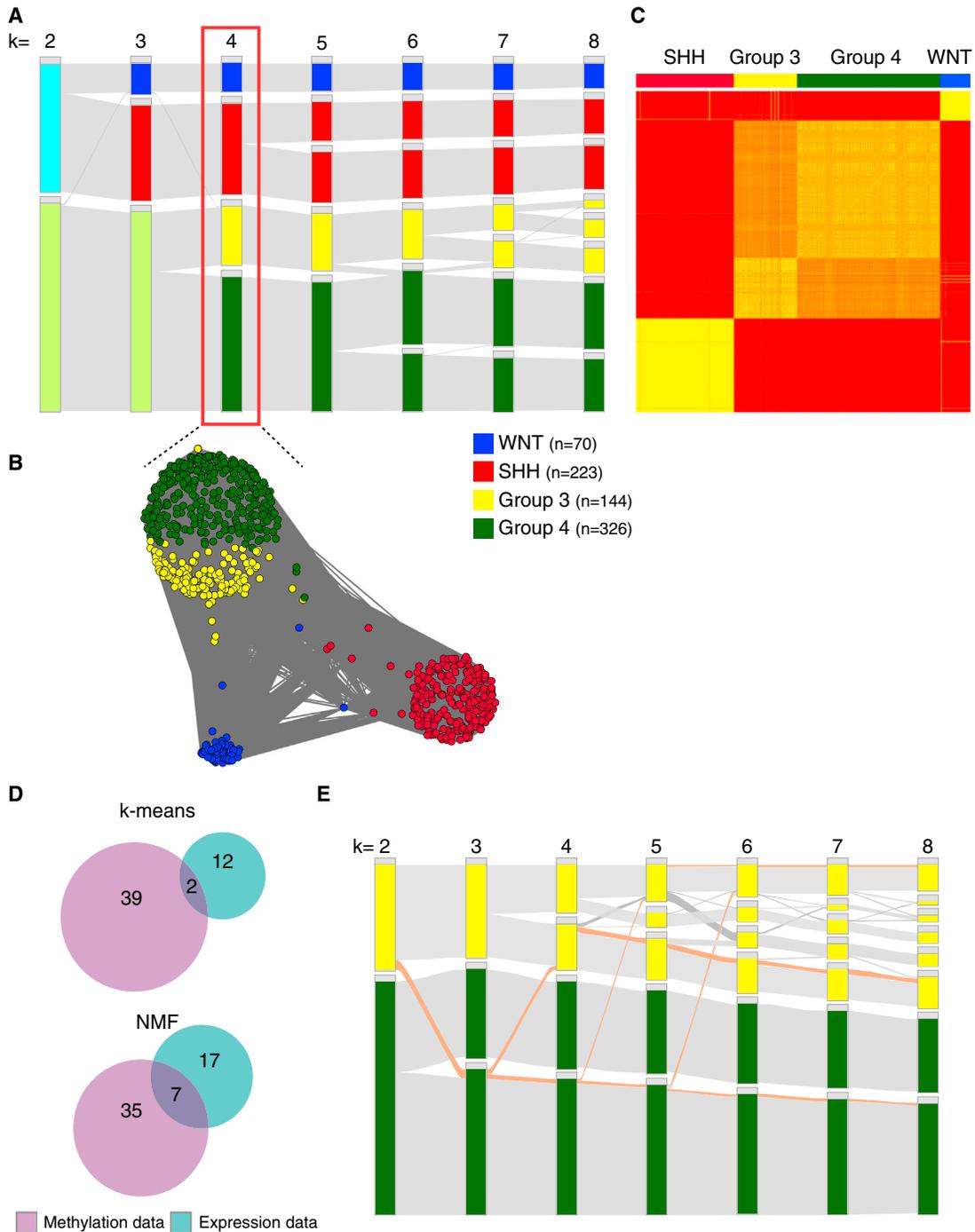


Figure 1. Clear Separation of the Four Medulloblastoma Subgroups through Integrative SNF Clustering

(A) Tumor clusters obtained by spectral clustering (for $k = 2$ to 8 groups) on the SNF network fused data obtained from both gene expression and DNA methylation data on 763 primary medulloblastomas. Relationships between tumors are indicated by the gray bars between columns. $k = 4$ (red box), defines the four recognized subgroups.

(B) Network representation of the relationships between tumors ($k = 4$). The shorter the edge between samples (nodes) is the more similar the samples are (only edges with a similarity value above the median value of all patient to patient similarity values are displayed).

(C) Heatmap representation of the sample-to-sample fused network data sorted by cluster for $k = 4$. Sample similarity is represented by red (less similar) to yellow (more similar) coloring inside the heatmap.

(D) Venn diagram showing the number of samples intermediate between groups 3 and 4 when using k-means or NMF clustering method on just expression or just methylation datasets of group 3 and 4 tumors ($n = 470$) between $k = 2$ and 3.

(legend continued on next page)

patient clusters. However, the subgroups identified using the two data types in isolation have not been compared head to head. More recently, methods of integrative clustering that analyze multiple data types in aggregate have been developed, including similarity network fusion (SNF) (Wang et al., 2014). Integrative approaches using multiple data types have been suggested to provide superior results compared with the analysis of single data types in isolation. SNF creates a unified view of patients based on multiple heterogeneous data sources, as it can integrate both gene- and non-gene-based data. SNF avoids the bias of genes or features pre-selection, is robust to different types of noise, is highly scalable, and has been shown to outperform other approaches for data integration (Wang et al., 2014).

Prior reports have recognized the existence of additional substructure within the four consensus subgroups, particularly within groups 3 and 4 (Cho et al., 2011). Consequently, a medulloblastoma consensus conference established that subdivisions within the known subgroups would be defined as subtypes, and labeled α , β , γ , δ , ϵ , etc. (Taylor et al., 2012). In this study our goal was to resolve intra-subgroup heterogeneity and identify biologically distinct and clinically relevant medulloblastoma subtypes by studying a very large cohort of primary tumor samples.

RESULTS

Integrated Clustering of Primary Medulloblastomas Recovers the Four Subgroups and Further Separates Group 3 from Group 4 Tumors

Through the Medulloblastoma Advanced Genomics International Consortium, we assembled a cohort of 763 primary frozen medulloblastoma samples with high-quality DNA and RNA, and generated genome-wide methylation and expression profiles. Of these, 491 had DNA copy-number profiles generated by Affymetrix SNP6 microarrays (Northcott et al., 2012b). Clinical data including age, tumor histology, metastatic status, and survival were available on 95.7%, 76.9%, 75.2%, and 82% of cases, respectively (Table S1). Arm-level somatic copy-number aberrations (SCNA) were inferred from methylation arrays in 100% of cases.

To these samples, we applied SNF to integrate both gene expression and DNA methylation data, followed by spectral clustering ranging from 2 to 12 groups. At $k = 4$, four distinct subgroups are clearly identified. Those groups correspond clinically and structurally to the previously described consensus subgroups: WNT ($n = 70$), SHH ($n = 223$), group 3 ($n = 144$), and group 4 ($n = 326$) (Figures 1A–1C and S1A–S1F) (Taylor et al., 2012).

Groups 3 and 4 are more similar to each other than to SHH and WNT (Figures 1B and 1C). We tested the stability of these core subgroups, by counting samples that switch subgroup affiliation when the number of clusters increases (Figure 1A). Following each sample from $k = 4$ to $k = 12$, no sample changed affiliation between WNT and SHH, while a small minority of samples moved between groups 3 and 4.

To determine the degree of overlap between groups 3 and 4, we undertook unsupervised clustering of 470 group 3 and 4 tumors using DNA methylation array data only, and then subsequently using transcriptional profiling data only. Both k -means and non-negative matrix factorization (NMF) consensus clustering revealed a small subset of tumors (2.9%–8.9%) that switched subgroup between $k = 2$ and $k = 3$ as determined through analysis of either transcriptional or methylation data (Figures S1G and S1H). Strikingly, the set of “ambiguous group 3–4 tumors” identified by gene expression profiling had very little overlap with those identified by DNA methylation profiling (Figure 1D) suggesting that the identification of the ambiguity may be a limitation of the particular type of measurement or data, rather than the identification of a truly distinct biological subtype. Examination of tumors within the “overlap” group does not reveal any demographic, clinical, or genetic commonalities, suggesting that it could be an artifact rather than a biologically discrete, clinically important group. Subsequent application of SNF and spectral clustering to this cohort of group 3 and 4 samples demonstrates that only 13/470 (2.8%) of samples change subgroup between $k = 2$ to $k = 3$, and of these 13 only 3 (0.64%) do not track back to their original subgroup when $k > 3$ (Figures 1E and S1I). We conclude that group 3 and group 4 medulloblastomas are stable, mostly non-overlapping molecular subgroups, and that SNF followed by spectral clustering is a more robust method of delineating subgroups than using a single data type in isolation.

Integrated Clustering Identifies 12 Medulloblastoma Subtypes

We applied SNF and spectral clustering within each of the four subgroups as defined by $k = 4$ across the entire cohort to determine the extent and nature of intra-subgroup heterogeneity. SNF and spectral clustering were selected to reduce the noise introduced by biased feature selection, and to leverage the full spectrum of our dataset. We identified clusters from $k = 2$ to $k = 8$ within each subgroup. In addition, we applied seven different machine-learning classifiers to predict the SNF subtypes. Cluster assignments from spectral clustering on the SNF fused similarity matrix was used as the “ground truth” subtype assignments. We split the dataset into a 70% training set and 30% testing set, trained the various classification models in 5-fold cross-validation on the training set and repeated the procedure 100 times (Table S2). We then applied the following criteria a priori to select the optimal number of subtypes: (1) how similar are the SNF clusters on the sample-to-sample heatmap? (2) How subtype specific are the broad and focal SCNA? (3) How relevant are the clinical associations? (4) How robustly can these subtypes be predicted using supervised machine learning? Using these criteria, we identified 12 subtypes: two WNT, four SHH, three group 3, and three group 4. For each solution, we identified focal SCNA from SNP6 data and arm-level copy-number gains and losses using copy-number states inferred from the methylation arrays.

(E) Tumor clusters obtained through spectral clustering on the SNF network fused data of group 3 and 4 samples ($n = 470$). A small minority of samples ($n = 13$, 2.8%) that were initially classified as group 3 samples at $k = 2$, subsequently move to group 4 at $k = 3$. Only 3/470 (0.64%) samples remain in group 4 after $k = 5$. These samples are tracked up to $k = 8$ (orange). See also Figures S1, S3 and Table S1.

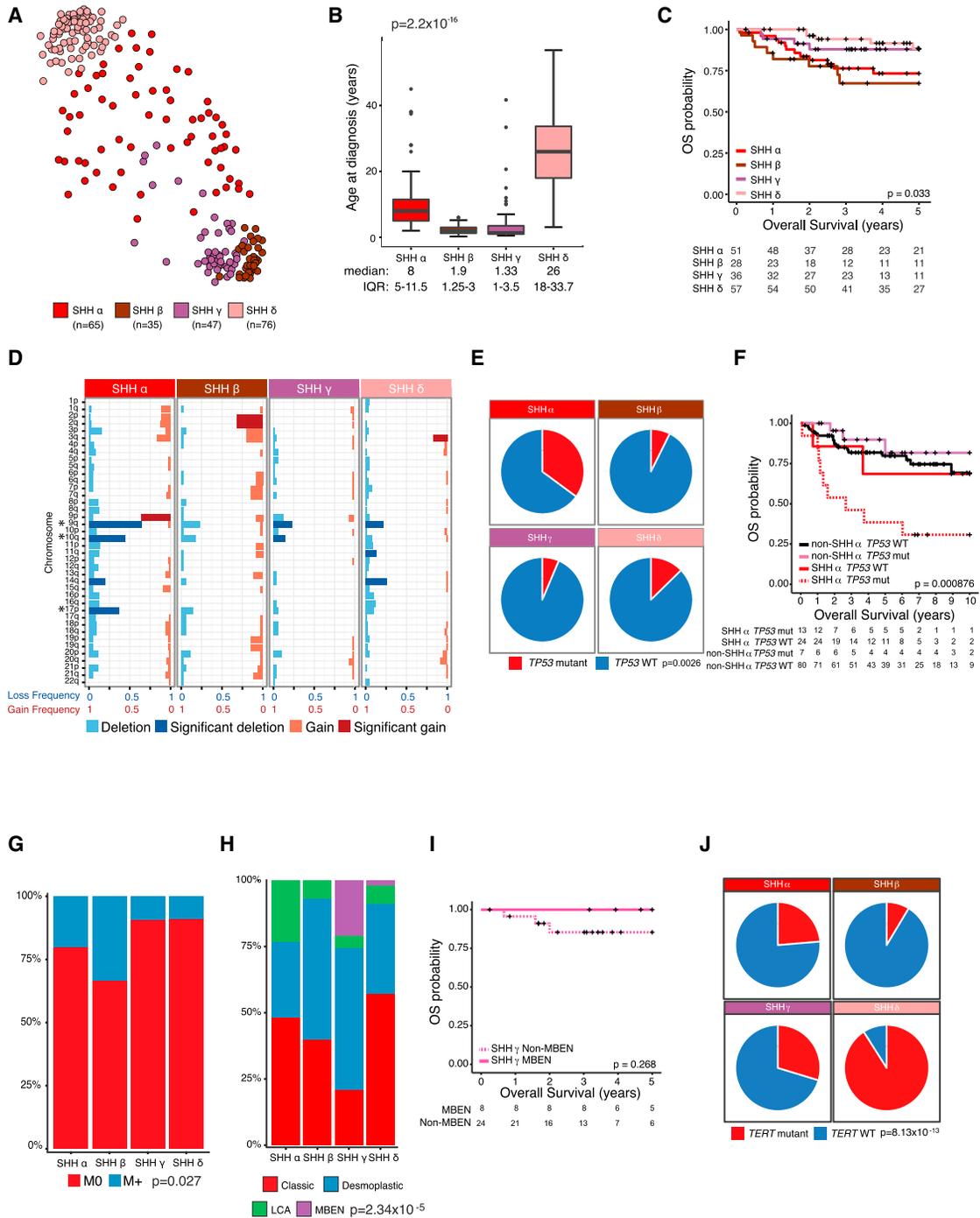


Figure 3. Clinical and Genomic Characteristics between Four SHH Medulloblastoma Subtypes

(A) Network representation map of $k = 4$ SNF-derived subtypes.
 (B) Age at diagnosis for SHH subtypes at $k = 4$ (Kruskal-Wallis test). Boxplot center lines show data median; box limits indicate the 25th and 75th percentiles; lower and upper whiskers extend 1.5 times the interquartile range (IQR) from the 25th and 75th percentiles, respectively. Outliers are represented by individual points.
 (C) Overall survival of SHH subtypes (log rank test). + indicates censored cases.
 (D) Frequency and significance of broad cytogenetic events across the four SHH subtypes. Darker bars show significant arm-level copy-number event ($q \leq 0.1$, chi-square test). * indicates key statistically significant arm gain or deletion.
 (E) Distribution of *TP53* mutations across SHH subtypes (Pearson's chi-square test).
 (F) Overall survival stratified by *TP53* mutation within SHH α and non-SHH α (log rank test). + indicates censored cases.
 (G) Incidence of metastatic dissemination at diagnosis across the four SHH subtypes (chi-square test).

(legend continued on next page)

Table S3). We evaluated the relationship between the associated genes and methylation probes in each subgroup. We first evaluated the number of associated genes that had methylation probes in their promoter region. Then we identified the subset of associated genes for which those probes were subgroup associated, and finally checked if we could detect an anti-correlation between the associated gene expression and the associated probe methylation levels. Only 3.7%, 8.3%, 6%, and 13% of *WNT*, *SHH*, group 3, and group 4 associated genes, respectively, follow all the criteria described above (Figure 2C). Therefore, only a small percentage of the associated genes are directly affected by DNA methylation. This is in support of both DNA methylation and gene expression contributing to the heterogeneity observed within each subgroup.

Integrative Clustering of DNA Methylation and Gene Expression Overcomes Discrepancies in Single Dataset Analysis at Defining Subtypes

To determine whether analysis of a single data type in isolation yielded similar results, we performed NMF clustering using gene expression or DNA methylation data individually. Using NMF clustering of the most variable expressed genes and methylated probes, we found that the two different types of data yield discordant subtypes as defined by both the cophenetic coefficient and silhouette value (>0.9) criteria (Figures S3A–S3D). In addition, the group memberships between the two modalities are divergent, indicating a lack of agreement between expression and methylation when analyzed in isolation (Figures S3E–S3H). When compared with the SNF subtypes, we found important differences, suggesting that both methylation and expression signatures contribute significantly and differently to define heterogeneity within the four subgroups; the data types provide distinct but complementary signals that improve over single-modality analyses. The subtypes identified by SNF are truly a combination of information present in both datasets, and therefore both data types are required to gauge the true intertumoral heterogeneity of medulloblastoma. For example, we observe that SHH α is mainly supported by the methylation data, but the defined group does not contain all SHH α samples (61%, Figure S3B). SHH δ is strongly supported by both the expression and methylation data (Figure S3B). In addition, groups 3 β and 3 γ are mainly defined by the signatures found in the expression data and do not separate well using the methylation data alone (Figure S3C). Finally, group 4 γ is very well supported by the methylation data, and corresponds to a group obtained with the expression data, but this latter group is missing 24.4% of group 4 γ samples (Figure S3D). Group 4 β is well supported by both data types (Figure S3D). We conclude that methylation and expression data are complementary, and an integrated approach allows a unified view of the underlying groups that is very valuable in elucidating heterogeneity within subgroups.

SHH Subtypes

Applying SNF and spectral clustering on SHH subgroup samples at $k = 4$ identified four clinically and cytogenetically distinct groups: SHH α ($n = 65$), SHH β ($n = 35$), SHH γ ($n = 47$), and SHH δ ($n = 76$) (Figures 3A, S4A, and S4B). SHH α tumors primarily affect children aged 3–16 years (Figure 3B), have the worst prognosis ($p = 0.03$, log rank test, Figure 3C), and are enriched for *MYCN* amplifications (SHH α 8/37, β 3/23, γ 0/29, δ 1/48; $p = 0.0034$ Pearson's chi-square test), and *GLI2* amplifications (SHH α 6/37, β 0/23, γ 0/29, δ 0/48; $p = 0.0002$ Pearson's chi-square test, Figure S4C; Table S4). Specific CNAs including 9q loss (SHH α 42/65, β 8/35, γ 11/47, δ 17/76; $p = 2.94 \times 10^{-7}$ Pearson's chi-square test), 10q loss (SHH α 29/65, β 6/35, γ 7/47, δ 6/76; $p = 1.54 \times 10^{-5}$ Pearson's chi-square test), 17p loss (SHH α 24/65, β 5/35, γ 3/47, δ 8/76; $p = 3.44 \times 10^{-5}$ Pearson's chi-square test, Figure 3D), and *YAP1* amplifications (SHH α 3/37, β 0/23, γ 0/29, δ 0/48; $p = 0.04$ Pearson's chi-square test, Figure S4C; Table S4) are also enriched in SHH α . The recent WHO classification includes SHH-activated *TP53* mutant tumors as a distinct category based on studies showing this group as being very high risk (Louis et al., 2016; Ramaswamy et al., 2016a; Zhukova et al., 2013). To further explore this association, *TP53* was sequenced across 145 SHH samples. *TP53* mutations are highly enriched in SHH α (SHH α 14/40, β 2/27, γ 2/31, δ 6/47; $p = 0.0026$ Pearson's chi-square test, Figure 3E; Table S5). When survival is analyzed stratified by *TP53* mutation and SHH α subtype, *TP53* mutations are only prognostic in SHH α (HR *TP53* mut versus WT: SHH α 6.006 [95% CI: 1.586–22.75; $p = 0.00832$] and non-SHH α 1.222 [95% CI: 0.2795–5.342; $p = 0.79$, Cox proportional hazards, Figure 3F]).

Interestingly, infant SHH tumors are mainly distributed across SHH β and SHH γ (age < 3: SHH α 5/65, β 23/35, γ 34/47, δ 0/76; $p = 2.2 \times 10^{-16}$ Pearson's chi-square test, Figure 3B), with disparate outcomes and copy-number profiles. SHH β tumors are frequently metastatic (33.3% versus 9.4% in SHH β and γ ; $p = 0.027$ Pearson's chi-square test, Figure 3G), harbor focal *PTEN* deletions (25% in SHH β versus none in γ), have multiple focal amplifications (Figure S4C; Table S4), and have a worse overall survival compared with SHH γ (HR of SHH β versus γ : 2.956 95% CI: 0.908–9.63; $p = 0.059$ Cox proportional hazards, Figure 3C). The difference in outcomes between SHH β and γ is possibly related to the increased rate of metastatic dissemination in SHH β , as there is a clear trend toward metastases being a marker of poor outcome within SHH β (HR of SHH β metastatic versus non-metastatic: 3.621 95% CI: 0.798–16.44; $p = 0.096$ Cox proportional hazards). Conversely, SHH γ have a relatively quiet copy-number landscape, with no recurrent amplifications, only one low-level recurrent focal deletion, and no significant arm-level gains (Figures 3D and S4C). Moreover, SHH γ are enriched for the MBEN (medulloblastoma with extensive nodularity) histology (20.9%; $p = 2.34 \times 10^{-5}$, Pearson's chi-square test, Figure 3H), which is known to portend more indolent clinical behavior (Rutkowski et al., 2010). Although almost all SHH tumors with MBEN histology ($n = 10$) were assigned to it,

(H) WHO histological classification at diagnosis across the four SHH subtypes (chi-square test).

(I) Overall survival within SHH γ stratified by MBEN histology (log rank test). + indicates censored cases.

(J) Distribution of *TERT* promoter mutations across SHH subtypes (Pearson's chi-square test).

See also Figures S4, S5; Tables S2, S4, and S5.

only a minority of SHH γ tumors have MBEN histology, demonstrating that histology alone is an inadequate surrogate to identify SHH γ tumors. The survival difference of SHH γ patients is not statistically significant between MBEN and non-MBEN tumors, suggesting that subtype affiliation is a more powerful biomarker than histopathology in infants with SHH medulloblastoma ($p = 0.268$, log rank test, Figure 3I). SHH δ are primarily composed of adults, have a favorable prognosis, and are strongly enriched for *TERT* promoter mutations (SHH α 6/34, β 2/22, γ 7/26, δ 38/42; $p = 8.13 \times 10^{-13}$, Pearson's chi-square test, Figure 3J).

To interrogate other possible solutions and to present the full results (Figures S5A–S5E), we also compared SHH subtypes when divided into three or five SNF groups. We refer to the clusters obtained by SNF for other numbers of groups ($k = 3$, $k = 5$ here) as c_1 , c_2 , c_3 , etc. (see Figures S4A and S4B). When comparing $k = 4$ with $k = 3$, SHH α and δ correspond closely to c_2 and c_1 , respectively, with c_3 representing a group of infants comprising SHH β and γ (Figures S4A, S4B, and S5A). SHH $k = 5$ reveals an additional group comprised primarily of a subset of SHH α patients with a group (c_3) enriched for 9q loss with a good prognosis and a second group (c_5) with a poor prognosis enriched for anaplasia (Figures S4B and S5C–S5E). Several machine-learning classifiers using both data types suggest poor confidence (<80%) in predicting the c_5 group. The machine-learning classifier with the best performance, elastic net (Zou and Hastie, 2005), is able to distinguish between four groups with >90% accuracy (Table S2). The identification of two groups of infant medulloblastoma with distinct clinical behavior allows for more precise and rational planning of clinical trials for infants with SHH medulloblastoma (Lafay-Cousin et al., 2016).

WNT Subtypes

We identify two WNT subtypes, WNT α ($n = 49$) and WNT β ($n = 21$) (Figures 4A, S6A, and S6B); WNT α is comprised mainly of children (Figure 4B), has similar survival as WNT β ($p = 0.5$, log rank test, Figure 4C), and has ubiquitous monosomy 6 (WNT α 48/49, β 6/21; $p = 2.365 \times 10^{-10}$ Pearson's chi-square test, Figure 4D). WNT β is enriched for older patients ($p = 4.013 \times 10^{-6}$, Kruskal-Wallis test, Figure 4B) who are frequently diploid for chromosome 6 (Figure 4D). Monosomy 6 has previously been described as a defining WNT medulloblastoma feature; clearly, patients with WNT β will be misdiagnosed if this criterion is used alone. Prior reports suggesting that adult WNT medulloblastoma might have a different biology and worse prognosis than childhood WNT medulloblastoma, are supported by our current analysis (Remke et al., 2011; Zhao et al., 2016). At $k = 3$, we observe a new group, comprised primarily of WNT β without monosomy 6 (Figures S6A and S6B); however, in the absence of any other defining feature or clear clinical relevance, we chose $k = 2$ as our preferred solution.

Group 3 Subtypes

Three very distinct subtypes of group 3 emerge from our analysis, each with characteristic copy-number and clinical variables: group 3α ($n = 67$), group 3β ($n = 37$), and group 3γ ($n = 40$) (Figures 5A, S6C, and S6D). A total of 60% of infants under the age of 3 years are in group 3α (age < 3: group 3α 14/63, 3β 4/36, 3γ 5/36; $p = 0.021$, Kruskal-Wallis test, Figure 5B).

Clinically, groups 3α and 3β have a more favorable prognosis compared with group 3γ (Figure 5C). Group 3β are slightly older ($p = 0.021$, Kruskal-Wallis test, Figure 5B), and are infrequently metastatic (group 3α 23/53, 3β 5/25, 3γ 15/30; $p = 0.058$ Pearson's chi-square test, Figure 5D). Group 3α and 3γ have a similar frequency of metastatic dissemination at diagnosis (Figure 5D). Chromosome 8q (*MYC* locus at 8q24) loss is more frequent in group 3α and gain more frequent in group 3γ (8q gain: group 3α 0/67, 3β 3/37, 3γ 22/40; $p = 2.2 \times 10^{-16}$ Pearson's chi-square test, Figure 5E), group 3β tumors have a higher frequency of activation of the *GFI1* and *GFI1B* oncogenes, previously shown to be drivers of group 3 through a process termed enhancer hijacking via focal gains and losses on chromosomes 1 and 9, with a paucity of arm-level chromosomal gains and losses (*GFI1* or *GFI1B* activation: group 3α 1/67, 3β 26/37, 3γ 3/40, $p < 2.2 \times 10^{-16}$ Pearson's chi-square test, Figures S7A and S7B) (Northcott et al., 2014). *OTX2* amplifications are also enriched in group 3β , as are losses of *DDX31* on chromosome 9; previously described to lead to activation of *GFI1B* through enhancer hijacking (*OTX2*: group 3α 0/35, 3β 6/28, 3γ 0/24; $p = 0.0013$; *DDX31* deletion: group 3α 1/35, 3β 9/28, 3γ 0/24; $p = 0.0031$ Pearson's chi-square test, Figure S7A; Table S4). Group 3γ have the worst prognosis ($p = 0.036$ log rank test, Figure 5C), a trend to enrichment of i17q (group 3α 17/67, 3β 5/37, 3γ 10/40; $p = 0.32$ Pearson's chi-square test, Figure 5E) and frequently harbor increased *MYC* copy number (group 3α 0/35, 3β 2/28, 3γ 5/24; $p = 0.012$, Figures 5F and S7A; Table S4), without other focal aberrations (Taylor et al., 2012). Group 3γ have a poor prognosis independent of *MYC* amplification, expanding the group of high-risk group 3 tumors beyond just *MYC* status ($p = 0.026$, log rank test, Figure 5G).

We find less support for other solutions of group 3, specifically $k = 2$ and $k = 4$ (Figures S6C, S6D, S7C, and S7D). At $k = 2$, we observe a group enriched for *MYC* amplification (c_1 0/38, c_2 7/48; $p = 0.014$ Pearson's chi-square test), and *GFI1* family of oncogene activations cluster together (*GFI1/1B* activation: c_1 1/71, c_2 29/73; $p = 1.14 \times 10^{-8}$ Pearson's chi-square test) without any meaningful clinical differences (Figure S7C). At $k = 4$, group 3α splits into two groups with minor contributions from the other two groups without any new meaningful clinical or copy-number enrichment (Figures S6D and S7D). In addition the elastic net classifier performs strongly at $k = 3$ (89%–98.8% per-group accuracy), while at $k = 4$ one group is less reliably predicted (72% accuracy, Table S2).

Group 4 Subtypes

Group 4 is the most prevalent subgroup comprising >40% of all medulloblastomas; previously described features include i17q, tandem duplications of *SNCAIP*, and high-level amplifications of *MYCN* and *CDK6* (Northcott et al., 2012b). We observe clear enrichment of key focal and arm-level SCNA at $k = 3$: group 4α ($n = 98$), group 4β ($n = 109$), and group 4γ ($n = 119$) (Figures 6A, S8A, and S8B). Clinically we observe group 4β have a slightly higher median age at diagnosis (8.22, 10, and 7 years for groups 4α , 4β , and 4γ ; $p = 1.34 \times 10^{-5}$ Pearson's chi-square test, Figure 6B); however, there is no statistically significant difference in the overall survival (Figure 6C) or rate of metastatic dissemination at diagnosis (groups 4α 30/75, 4β 35/86, 4γ 36/94; $p = 0.94$ Pearson's chi-square test, Figure 6D). Group 4α are

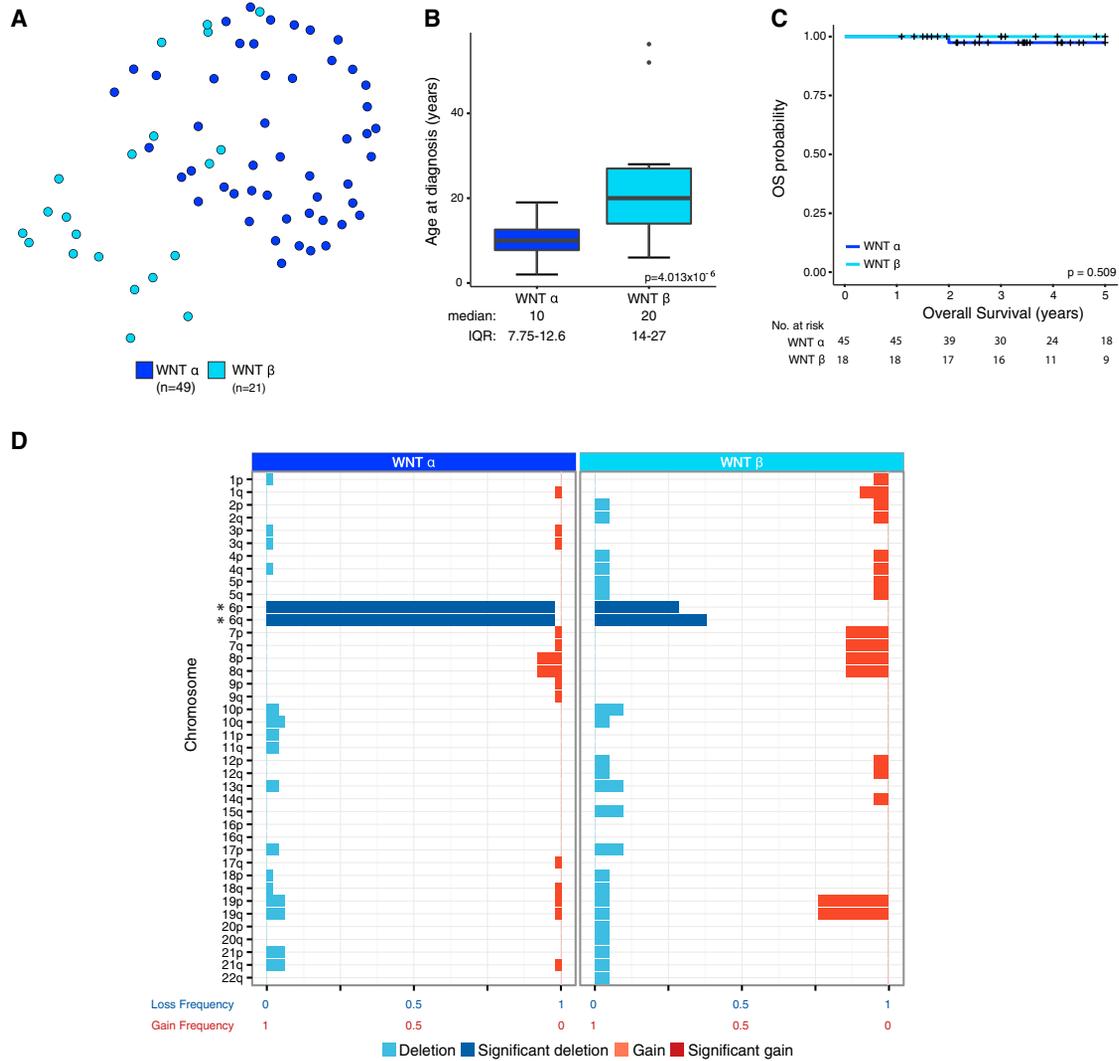


Figure 4. Clinical and Genomic Characteristics between Two WNT Medulloblastoma Subtypes

(A) Network representation map of $k = 2$ SNF-derived subtypes.

(B) Age at diagnosis for WNT subtypes at $k = 2$ (Mann-Whitney U test). Boxplot center lines show data median; box limits indicate the 25th and 75th percentiles; lower and upper whiskers extend 1.5 times the interquartile range (IQR) from the 25th and 75th percentiles, respectively. Outliers are represented by individual points.

(C) Overall survival comparing WNT α with WNT β (log rank test). + indicates censored cases.

(D) Frequency and significance of broad cytogenetic events across the two WNT subtypes. Darker bars show significant arm-level copy-number events ($q \leq 0.1$, chi-square test). * indicates key statistically significant arm gain or deletion.

See also Figure S6.

enriched for *MYCN* amplifications (11/66, compared with none in group 4 β and 4 γ ; $p = 2.46 \times 10^{-6}$ Pearson's chi-square test, Figure S8C; Table S4). Group 4 α and 4 γ are strongly enriched for 8p loss (group 4 α 47/98, 4 β 24/109, 4 γ 87/119; $p = 1.22 \times 10^{-13}$ Pearson's chi-square test) and 7q gain (group 4 α 57/98, 4 β 9/109, 4 γ 62/119; $p = 9.5 \times 10^{-31}$, Pearson's chi-square test, Figure 6E). Group 4 β are strongly enriched for *SNCAIP* duplications (group 4 α 4/66, 4 β 11/74, 4 γ 0/73; $p = 0.0019$ Pearson's chi-square test) and almost ubiquitous i17q (group 4 α 40/98, 4 β 87/109, 4 γ 31/119; $p = 9.75 \times 10^{-16}$ Pearson's chi-square test) with a paucity of other SCNA (Figures 6E and S8C; Table S4). In addition, groups 4 α and 4 γ are enriched for focal *CDK6*

amplifications (group 4 α 4/66, 4 β 0/74, 4 γ 6/73; $p = 0.051$ Pearson's chi-square test, Figure S8C; Table S4). Previous studies have suggested *GFI1* and *GFI1B* activation to be present in group 4, however we see GFI activation to be largely restricted to group 3 β (Figure S8D).

At $k = 2$, we observe groups 4 α and 4 γ forming one group, and group 4 β being largely preserved (Figures S8A, S8B, and S8E). At $k = 4$, group 4 β continues to segregate from the other groups; however, no new groups emerge with any significant clinical or copy-number differences (Figures S8A, S8B, and S8F). Due to the enrichment of key SCNA at $k = 3$, we chose this as our preferred solution. Moreover, our classifier exhibits a decline in

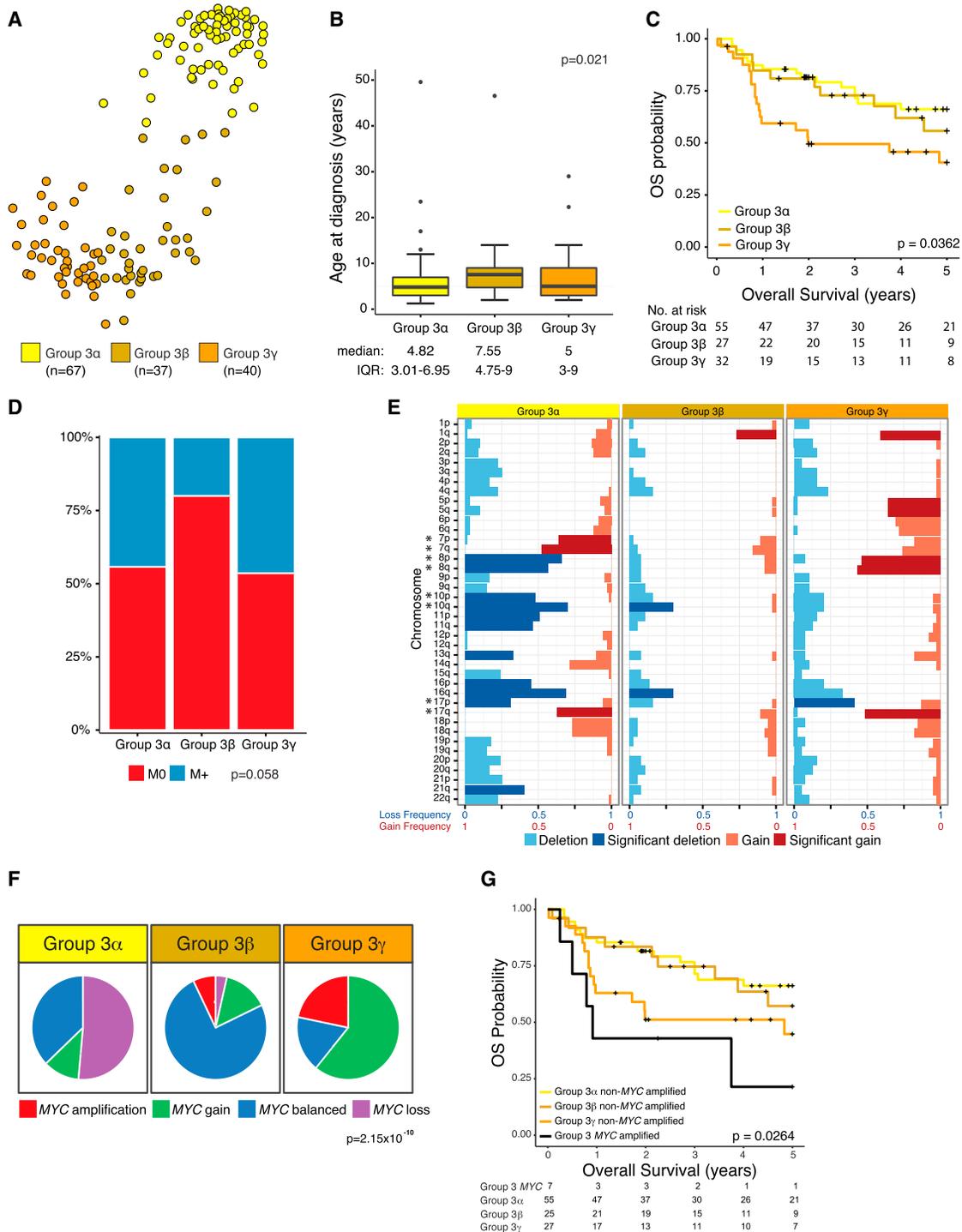


Figure 5. Clinical and Genomic Characteristics between Three Group 3 Medulloblastoma Subtypes

(A) Network representation map of $k=3$ SNF-derived subtypes.

(B) Age at diagnosis of group 3 subtypes at $k=3$ (Kruskal-Wallis test). Boxplot center lines show data median; box limits indicate the 25th and 75th percentiles; lower and upper whiskers extend 1.5 times the interquartile range (IQR) from the 25th and 75th percentiles, respectively. Outliers are represented by individual points.

(C) Overall survival of group 3 subtypes (log rank test). + indicates censored cases.

(D) Incidence of metastatic dissemination at diagnosis for the three group 3 subtypes (chi-square test).

(E) Frequency and significance of broad cytogenetic events across the group 3 subtypes. Darker bars show significant arm-level events ($q \leq 0.1$, chi-square test).

* indicates key statistically significant arm gain or deletion.

(legend continued on next page)

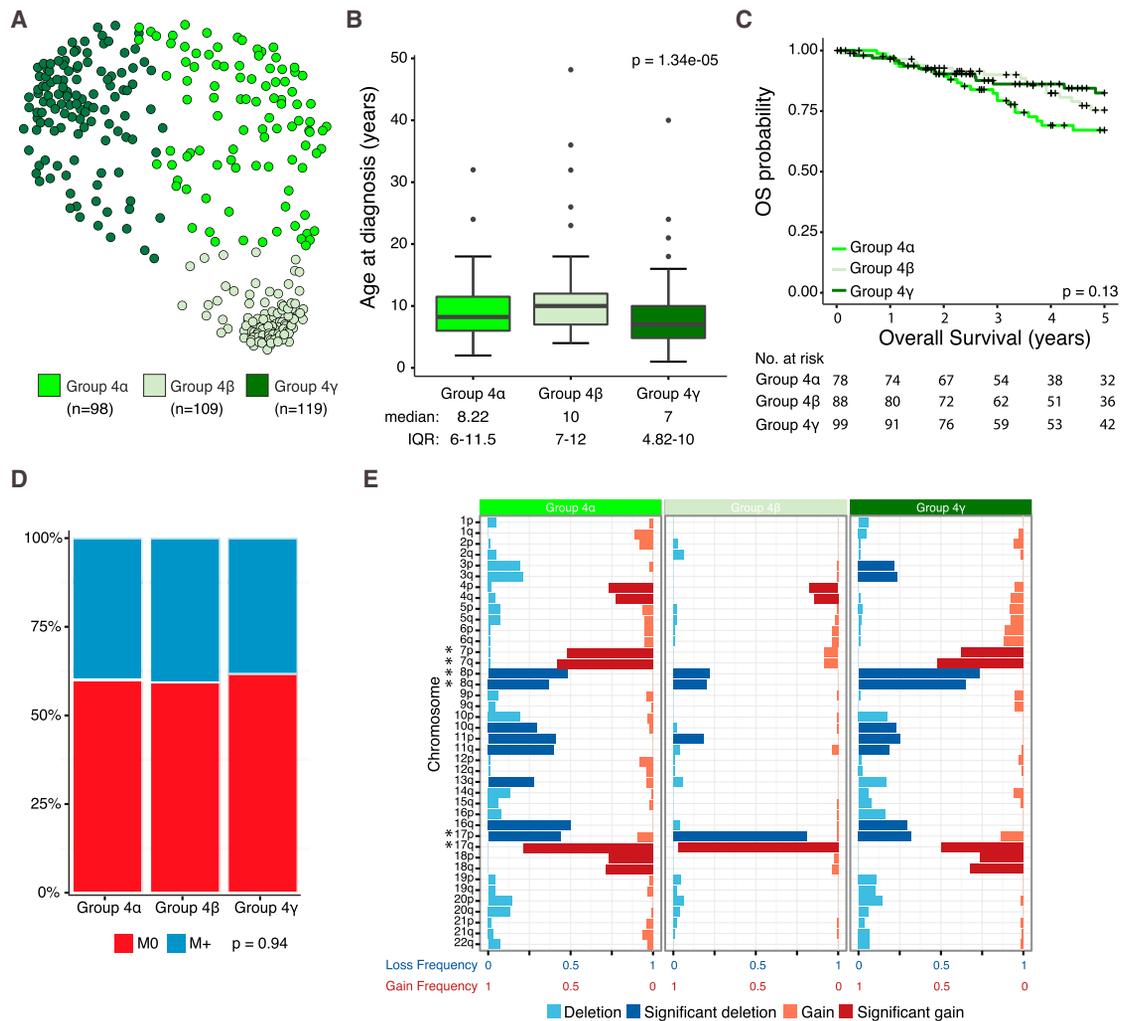


Figure 6. Clinical and Genomic Characteristics of the Three Group 4 Medulloblastoma Subtypes

(A) Network representation map of $k = 3$ SNF-derived subtypes.
 (B) Age at diagnosis of group 4 subtypes at $k = 3$ (Kruskal-Wallis test). Boxplot center lines show data median; box limits indicate the 25th and 75th percentiles; lower and upper whiskers extend 1.5 times the interquartile range (IQR) from the 25th and 75th percentiles, respectively. Outliers are represented by individual points.
 (C) Overall survival of group 4 subtypes (log rank test). + indicates censored cases.
 (D) Incidence of metastatic dissemination at diagnosis across the three group 4 subtypes (chi-square test).
 (E) Frequency and significance of broad cytogenetic events across the three group 4 subtypes. * indicates key statistically significant arm gain or deletion. Darker bars show significant arm-level events ($q \leq 0.1$, chi-square test).
 See also [Figure S8](#) and [Tables S2, S4](#).

confidence at $k = 4$, suggesting these groups are not as robust as $k = 3$ ([Table S2](#)).

Comparable Subtypes with Key Clinical Differences Are Identified by Other Integrative Analyses

Two other integrative clustering methods have been employed by the The Cancer Genome Atlas (TCGA) consortium in previous studies of other cancer histologies. We applied both methods to

our dataset; when applying the cluster of clusters (COCA) method used by TCGA in low-grade glioma and pan-cancer studies ([Bret et al., 2015](#); [Hoadley et al., 2014](#)) we observed that the method was quite limited in the potential to leverage information from our two data types in the current manuscript. The COCA subgroups were driven by the samples that agree or disagree between the two data types clustered in isolation, which is the COCA input. COCA failed to identify one SHH infant subtype or group 3β.

(F) Distribution of MYC amplifications across group 3 subtypes (Pearson’s chi-square test).
 (G) Overall survival of group 3 subtypes without MYC amplifications for each subtype compared with MYC-amplified tumors (log rank test). + indicates censored cases.
 See also [Figures S6](#) and [S7](#); [Tables S2](#) and [S4](#).

iCluster was used successfully by TCGA to identify relevant subtypes (Collisson et al., 2014). When applying iCluster to our dataset, at $k = 4$, the four groups did not have the demographics and SCNA consistent with the four previously described groups. When comparing the four iCluster groups with those defined by SNF, WNT and group 3 do not separate, and SHH comprises two groups. When we analyze the iCluster results for five groups, we recover two SHH groups, plus WNT, group 3, and group 4, which in this case corresponds very well to the SNF subgroup (when considering the two SHH groups together). We then asked if we could recover similar subtypes to SNF using iCluster. As we could not recover the four main groups, subgroups defined by SNF were then individually analyzed using iCluster. We observe a near 80% concordance with the SNF subtypes. The childhood and the adult SHH subtypes as well as the group 4 subtypes are recapitulated (along with a single SHH infant group). However, we identified key differences particularly within the WNT, SHH, and group 3 subgroups. Only one WNT group is identified, the two infant SHH subtypes are not identified, and the two distinct group 3 subtypes with *MYC* amplifications and *GFI1* activation are not observed. Clearly, the SNF method is superior at leveraging information of multiple datasets to identify meaningful groups of patients in a cancer cohort, specifically in a medulloblastoma cohort.

Differential Pathway Activation Defines Subtypes across All Four Medulloblastoma Subgroups

Pathway enrichment analysis was performed for each of the identified subtypes across all four subgroups using the top 10% of associated genes across each subtype. We observe several significantly enriched pathways for all identified subtypes (adj. p value < 0.05), supporting subtype-specific biological processes and transcriptional networks (Figures 7A–7D). In particular, in SHH we observe several pathways enriched in SHH β and γ , with developmental pathways more enriched in SHH γ over β (Figure 7A). Genes involved in DNA repair and cell cycle are significantly enriched in SHH α . Several actionable pathways, as defined by the availability of approved drugs, are subtype specific. Specifically, sumoylation is enriched in SHH α , ion channels are enriched in SHH β and γ , and telomere maintenance is enriched in SHH α and δ . Receptor tyrosine kinase signaling is enriched in SHH γ and, to a lesser extent, in β . DNA repair pathways are enriched in SHH α , suggesting that strategies to inhibit the DNA damage response and increase replicative stress are more likely to be effective in this group.

Group 3 α tumors are enriched for photoreceptor, muscle contraction, and primary cilium-related genes (Figure 7B). Pathways involved in protein translation are enriched in groups 3 β and 3 γ , which are potentially actionable using modulators of protein synthesis such as proteasome inhibitors. Telomere maintenance is also more enriched in group 3 γ , suggesting that telomerase inhibition may only be effective in one group. Several pathways are identified across group 4 subtypes, which, coupled with subtype-specific copy-number enrichment, further supports the existence of three group 4 subtypes (Figure 7C). Actionable pathways restricted to particular subtypes include MAPK and FGFR1 signaling in group 4 β and PI3K-AKT signaling and ERBB4-mediated nuclear signaling in group 4 γ . Cell migration pathways are more enriched in group 4 α .

DISCUSSION

Our study identifies and delineates the intertumoral heterogeneity present within medulloblastoma subgroups. Leveraging a large cohort of medulloblastomas profiled by combined gene expression and DNA methylation, we have identified different subtypes within each of the four core subgroups. These subtypes have particular clinical and copy-number features, which allow for a refinement in our understanding of the genomic landscape of medulloblastoma (Figure 8). Combining expression and methylation data using SNF adds further proof that groups 3 and 4 are largely different biological entities. The deeper we go in clustering medulloblastoma samples, the less consistent the groups become. This is exemplified by poor predictability of putative subtypes when a large number of subtypes is assumed. Defining clinical features and CNAs also tend to lose their distinctive profiles as we increase the number of clusters, suggesting that heterogeneity is bounded by a discrete number of optimal groups.

Comparison of SNF with consensus clustering of either gene expression or DNA methylation data analyzed in isolation clearly suggests that an integrated approach provides a much more refined and accurate classification. This is particularly striking when evaluating the boundary between groups 3 and 4, where samples that are deemed indeterminate using gene expression and DNA methylation in isolation are largely non-overlapping. Moreover, in elucidating the heterogeneity within subgroups, we observe significant disagreement between gene expression and DNA methylation in isolation, suggesting that each data type makes a unique and non-redundant contribution to defining the subtypes. The very low number of samples that change subgroup affiliation using SNF strongly advocates that definition of these two groups is largely enhanced using an integrative approach. A limitation of our approach is the bulk analysis of samples. At a subclonal level, a greater degree of overlap across groups 3 and 4 cannot be discounted. More detailed analysis at a cellular level, specifically applying single-cell methods, will help delineate the full subclonal structure, potentially uncovering subsets of group 3 and 4 samples with common mechanisms and cellular origins. Further studies integrating emerging technologies such as long non-coding RNA, proteomics, and histone modifications may allow an even more refined description of the medulloblastoma landscape; however, the large cohorts of frozen tissue required for these studies are presently not available.

The identification of subtypes has significant biological and clinical implications. Several previously described copy-number alterations within medulloblastoma subgroups such as amplifications/gains of *MYC*, *MYCN*, *OTX2*, *CDK6*, *SNCAIP*, and *ACVR1*, as well as several arm-level events including *i17q* clearly segregate between subtypes (Northcott et al., 2012b). Our identification of unique cytogenetic aberrations that occur in concert, as well as specific biological pathways enriched within specific subtypes, will serve to inform creation of rational preclinical models that closely mirror the human diseases. Several of these aberrations are actionable and largely restricted to subtypes, which will also allow for a more personalized treatment approach. Several subtypes, particularly in SHH and group 3, have clear and drastic clinical and prognostic differences, which

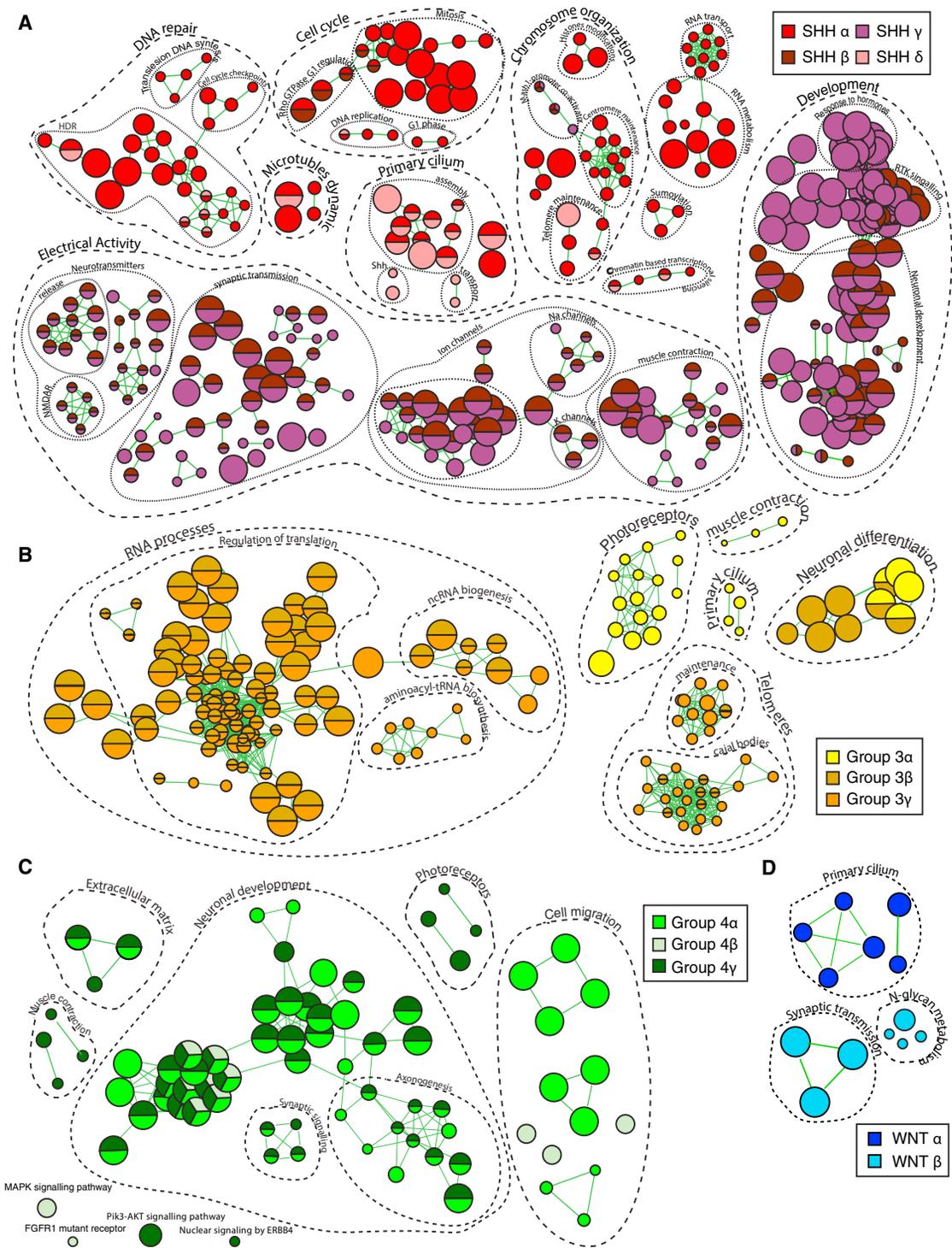


Figure 7. Subtype-Enriched Pathways

(A–D) Enrichment maps representing biological processes and pathways enriched in subtype-specific upregulated genes for SHH subtypes (A), group 3 subtypes (B), group 4 subtypes (C), and WNT subtypes (D). Each node represents a process or pathway; nodes with many shared genes are grouped and labeled by biological theme. Processes and pathways connected at edges have genes in common. Nodes are colored according to the subtype(s) in which the process is enriched; processes enriched in more than one subtype have multiple colors. Nodes sizes are proportional to the number of genes in each process, in each subgroup. Enriched processes were determined with g:Profiler (FDR-corrected q value < 0.05) and visualized with the Enrichment Map app in Cytoscape. Connected nodes and unconnected but actionable nodes are shown.

Subgroup		WNT		SHH				Group 3			Group 4		
Subtype		WNT α	WNT β	SHH α	SHH β	SHH γ	SHH δ	Group 3 α	Group 3 β	Group 3 γ	Group 4 α	Group 4 β	Group 4 γ
Subtype proportion													
Subtype relationship													
Clinical data	Age												
	Histology			LCA Desmoplastic	Desmoplastic	MBEN Desmoplastic	Desmoplastic						
	Metastases	8.6%	21.4%	20%	33%	8.9%	9.4%	43.4%	20%	39.4%	40%	40.7%	38.7%
	Survival at 5 years	97%	100%	69.8%	67.3%	88%	88.5%	66.2%	55.8%	41.9%	66.8%	75.4%	82.5%
Copy number	Broad	6 ⁻		9q ⁻ , 10q ⁻ , 17p ⁻		Balanced genome		7 ⁺ , 8 ⁻ , 10 ⁻ , 11 ⁺ , i17q			7q ⁺ , 8p ⁻ , i17q		
	Focal			MYCN amp, GLI2 amp, YAP1 amp		PTEN loss		10q22 ⁻ , 11q23.3 ⁻			OTX2 gain, DDX31 loss		
Other events				TP53 mutations				TERT promoter mutations			High GF11/B expression		

Age (years): 0-3 >3-10 >10-17 >17

Figure 8. Graphical Summary of the 12 Medulloblastoma Subtypes

Schematic representation of key clinical data, copy-number events, and relationship between the subtypes inside each of the four medulloblastoma subgroups. The percentages of patients presenting with metastases and the 5-year survival percentages are presented. The age groups are: infant 0–3 years, child >3–10 years, adolescent >10–17 years, and adult >17 years.

will allow for more robust risk stratification in future clinical trials. Furthermore, a major hurdle to clinical trial design has been the overlap of groups 3 and 4 in current studies, which if applied today would make strata assignment difficult. The next generation of clinical trials for high-risk medulloblastoma will involve subgroup-specific therapies. The inability to stratify 10% of patients to either groups 3 or 4 has the potential to either deprive a patient of an innovative therapy or, of more concern, expose a child to an inappropriate escalation or de-escalation of therapy.

Clinically, our observed groups have immediate implications. It has been shown that *TP53* mutations are highly prognostic in SHH. We extend these findings whereby *TP53* mutations are not only enriched in SHH α but also only prognostic in SHH α . This is highly relevant for clinical trial design, where *TP53* mutant SHH has been identified as a very-high-risk group to be prioritized for novel therapies in both Europe and North America (Ramaswamy et al., 2016a); clearly, the observation that *TP53* mutations are highly enriched and prognostic in SHH α has significant implications. A limitation of this is the absence of germline status, which, based on previous studies, are likely *TP53* mutant enriched in SHH α .

The identification of two infant SHH groups has clear and immediate clinical significance. Currently, infant medulloblastomas are stratified by the presence or absence of desmoplastic morphology. However, several reports have suggested that in-

fant SHH as a whole have a favorable prognosis independent of morphology. Our results suggest that clinical risk stratification can be refined by incorporation of integrated subtypes, whereby SHH γ are clearly a very low-risk group and could be spared the toxic effects of high-dose chemotherapy. Our observation that MBEN histology is almost exclusive to SHH γ , but represent a minority of cases within SHH γ , has significant implications for clinical trials. Current infant clinical trials stratify patients based on either classic or desmoplastic/MBEN histology. Indeed, the frequency of desmoplastic histology is similar across all four SHH subtypes, despite significant differences in survival between SHH subtypes. The most recent infant medulloblastoma study from the Children's Oncology Group ACNS1221 (NCT02017964) was closed prematurely due to an excess of relapses. This study selected infants with a "desmoplastic" morphology for treatment de-escalation, of which the vast majority are SHH. Indeed, our identification of two infant subtypes of SHH represents an example where more robust risk stratification has the potential to accurately select patients for de-escalation of therapy in future clinical trials. Overall, this further supports the idea that the incorporation of molecular stratification rather than subjective morphology alone has the potential for immediate clinical benefit.

Similarly, for group 3, we identify a high-risk group that is enriched for *MYC* amplification, but for which not all patients are *MYC* amplified. Interestingly, the majority of in vitro cell lines of

medulloblastoma do not represent the clear intertumoral heterogeneity, but rather are *MYC*-amplified or *MYC*-activated models that actually represent only group 3 γ . The identification of significant heterogeneity across group 3 underlies the urgent need to develop preclinical models that faithfully recapitulate the different subtypes within each subgroup. In group 4, there are currently no robust preclinical models, and the subgroups we describe, specifically the mutually exclusivity of *MYCN* amplifications and *SNCAIP* duplications, may help with future modeling.

Taken together, our results highlight the power of combining multiple data types compared with the use of single data types in isolation. This approach has identified that there may be a limit to the degree of substructure across medulloblastoma samples; however, only a study with a much larger cohort could fully assess the extent of intertumoral heterogeneity within the subgroups. We identify clinically important substructure within subgroups, which will allow further refinement of our biological and clinical risk stratification schemes. The identification of homogeneous subtypes may simplify the identification of targets for therapy, and could allow for therapies effective across subtypes. The development of reliable biomarkers to identify subtypes will provide much needed prognostic information for patient stratification, particularly in regard to SHH and group 3 medulloblastoma.

STAR★METHODS

Detailed methods are provided in the online version of this paper and include the following:

- **KEY RESOURCES TABLE**
- **CONTACT FOR REAGENT AND RESOURCE SHARING**
- **EXPERIMENTAL MODEL AND SUBJECT DETAILS**
 - Acquisition of Patient Samples
- **METHOD DETAILS**
 - Nucleic Acid Extraction
 - Expression and Methylation Data
 - *TERT* Promoter and *TP53* Sequencing
- **QUANTIFICATION AND STATISTICAL ANALYSIS**
 - Microarray Gene Expression Analysis
 - Genome Wide Methylation Analysis
 - Methylation Array Copy Number Analysis
 - SNP6 Copy Number Analysis
 - Clinical Correlation and Survival Analysis
 - Group 3 and Group 4 Analysis
 - Similarity Network Fusion Analysis (SNF)
 - Groups Visualization Using Stratomex
 - Network Visualization with Cytoscape
 - Relationship between Associated Genes and Probes
 - Pathway Enrichment Analysis
 - Classifier Description
 - Training and Selection of the Classifiers
 - COCA Analysis
 - iCluster Analysis
- **DATA AND SOFTWARE AVAILABILITY**

SUPPLEMENTAL INFORMATION

Supplemental Information includes eight figures and five tables and can be found with this article online at <http://dx.doi.org/10.1016/j.ccell.2017.05.005>.

AUTHOR CONTRIBUTIONS

Conceptualization, F.M.G.C., M.R., V.R., and M.D.T.; Methodology, F.M.G.C., M.R., L.R., D.J.H.S., B.L., J.T., A.S.Mor., S.Ag., E.B., and V.R.; Investigation, F.M.G.C., M.R., L.R., J.P., L.G., C.N., A.S.Mor., Y.Y.T., C.M.K., H.F., K.I., and J.R.; Data Curation, F.M.G.C., B.L., A.S.Mor., H.F., and V.R.; Formal Analysis: F.M.G.C., L.R., D.J.H.S., L.G., J.R., and V.R.; Validation, F.M.G.C., L.R., V.R., and A.G.; Writing – Original Draft, F.M.G.C., V.R., and M.D.T.; Writing – Review & Editing, F.M.G.C., C.D., J.R., A.G., V.R., and M.D.T.; Resources, B.-K.C., S.-K.K., K.-C.W., J.Y.L., W.A.G., M.P.-P., A.V., C.F.-C., A.J., C.G., A.A.N.R., K.K.W.L., H.-K.N., C.G.E., I.F.P., R.L.H., G.Y.G., J.M.O., S.L., W.A.W., B.L., L.B.C., R.C.T., M.K.C., R.V., P.H., M.-L.C.v.V., J.M.K., P.J.F., Y.S.R., T.K., E.L.-A., K.Z., J.S., G.F., M.M., E.G.v.M., S.O., T.S., A.K., M.Z., J.R.L., J.B.R., N.J., S.A.I., J.M., T.E.v.M., S.J., A.S.Moo., A.R.H., J.A.C., D.P.C.T., C.G.C., M.F., J.P., C.C.F., A.G.S., L.M., L.M.L., H.W., H.N., S.K.E., M.P.-D., F.C.P.d.L., S.R., M.Z., A.L., A.H., C.E.H., U.T., E.B., U.B., P.B.D., and J.T.R.; Project Administration, A.G., V.R., and M.D.T.; Supervision, G.D.B., A.G., V.R., and M.D.T.; Funding Acquisition, V.R. and M.D.T.

ACKNOWLEDGMENTS

M.D.T. is supported by funds from the Garron Family Chair in Childhood Cancer Research at The Hospital for Sick Children and The University of Toronto, and operating funds from the NIH (R01CA159859 and R01CA148699), The Terry Fox Research Institute, The Brain Tumor Foundation of Canada, The McLaughlin Center, Worldwide Cancer Research, The Canadian Institutes of Health Research, and the Pediatric Brain Tumor Foundation. M.D.T. is also supported by a Stand Up To Cancer - St. Baldrick's Pediatric Dream Team Translational Research Grant (SU2C-AACR-DT1113). Stand Up To Cancer is a program of the Entertainment Industry Foundation administered by the American Association for Cancer Research. V.R. is supported by Meagan's Walk, an Alex's Lemonade Stand Young Investigator Award, a Garron Family Cancer Center Pitblado Discovery Grant, and a Collaborative Ependymoma Research Network basic science fellowship. V.R. and M.D.T. are supported by the Swifty Foundation. F.M.G.C. is supported by the Stephen Buttrum Brain Tumor Research Fellowship, granted by Brain Tumor Foundation of Canada. M.R. is supported by a fellowship from the Mildred Scheel Cancer Foundation and operating funds from the Pediatric Brain Tumor Foundation. J.R. and K.I. were supported by NSERC Discovery Grant RGPIN-2016-06485. J.R. was supported by Operating Grant 21089 of the Cancer Research Society. A.K. was supported by the Hungarian Brain Research Program (grant no. KTIA_13_NAP-A-V/3), the TÁMOP-4.2.2.A-11/1/KONV-2012-0025 project, and the János Bolyai Scholarship of the Hungarian Academy of Sciences. K.Z. acknowledges research support from the project OPVK CZ.1.07/2.3.00/20.0183. E.V.M. is funded by St. Baldrick's Foundation, and NIH R01 NS084063. We thank Susan Archer for technical writing.

Received: September 23, 2016

Revised: March 24, 2017

Accepted: May 8, 2017

Published: June 12, 2017

REFERENCES

- Aryee, M.J., Jaffe, A.E., Corrada-Bravo, H., Ladd-Acosta, C., Feinberg, A.P., Hansen, K.D., and Irizarry, R.A. (2014). Minfi: a flexible and comprehensive Bioconductor package for the analysis of Infinium DNA Methylation microarrays. *Bioinformatics* 30, 1363–1369.
- Brat, D.J., Verhaak, R.G., Aldape, K.D., Yung, W.K., Salama, S.R., Cooper, L.A., Rheinbay, E., Miller, C.R., Vitucci, M., Morozova, O., et al. (2015). Comprehensive, integrative genomic analysis of diffuse lower-grade gliomas. *N. Engl. J. Med.* 372, 2481–2498.
- Cho, Y.-J., Tsherniak, A., Tamayo, P., Santagata, S., Ligon, A., Greulich, H., Berhoukim, R., Amani, V., Goumnerova, L., Eberhart, C.G., et al. (2011). Integrative genomic analysis of medulloblastoma identifies a molecular subgroup that drives poor clinical outcome. *J. Clin. Oncol.* 29, 1424–1430.

- Collisson, E.A., Campbell, J.D., Brooks, A.N., Berger, A.H., Lee, W., Chmielecki, J., Beer, D.G., Cope, L., Creighton, C.J., Danilova, L., et al. (2014). Comprehensive molecular profiling of lung adenocarcinoma. *Nature* **511**, 543–550.
- Dai, M., Wang, P., Boyd, A.D., Kostov, G., Athey, B., Jones, E.G., Bunney, W.E., Myers, R.M., Speed, T.P., Akil, H., et al. (2005). Evolving gene/transcript definitions significantly alter the interpretation of genechip data. *Nucleic Acids Res.* **33**, e175.
- Gaujoux, R., and Seoighe, C. (2010). A flexible R package for nonnegative matrix factorization. *BMC Bioinformatics* **11**, 367.
- Gautier, L., Cope, L., Bolstad, B.M., and Irizarry, R.A. (2004). *affy* — analysis of Affymetrix GeneChip data at the probe level. *Bioinformatics* **20**, 307–315.
- Gevaert, O. (2015). *MethylMix*: an R package for identifying DNA methylation-driven genes. *Bioinformatics* **31**, 1839–1841.
- Hovestadt, V., Remke, M., Kool, M., Pietsch, T., Northcott, P.A., Fischer, R., Cavalli, F.M., Ramaswamy, V., Zapatka, M., Reifemberger, G., et al. (2013). Robust molecular subgrouping and copy-number profiling of medulloblastoma from small amounts of archival tumour material using high-density DNA methylation arrays. *Acta Neuropathol.* **125**, 913–916.
- Hoadley, K.A., Yau, C., Wolf, D.M., Cherniack, A.D., Tamborero, D., Ng, S., Leiserson, M.D., Niu, B., McLellan, M.D., Uzunangelov, V., et al. (2014). Multiplatform analysis of 12 cancer types reveals molecular classification within and across tissues of origin. *Cell* **158**, 929–944.
- Kauffmann, A., Gentleman, R., and Huber, W. (2009). *arrayQualityMetrics* — a bioconductor package for quality assessment of microarray data. *Bioinformatics* **25**, 415–416.
- Kool, M., Jones, D.T., Jager, N., Northcott, P.A., Pugh, T.J., Hovestadt, V., Piro, R.M., Esparza, L.A., Markant, S.L., Remke, M., et al. (2014). Genome sequencing of SHH medulloblastoma predicts genotype-related response to smoothened inhibition. *Cancer Cell* **25**, 393–405.
- Lafay-Cousin, L., Smith, A., Chi, S.N., Wells, E., Madden, J., Margol, A., Ramaswamy, V., Finlay, J., Taylor, M.D., Dhall, G., et al. (2016). Clinical, pathological, and molecular characterization of infant medulloblastomas treated with sequential high-dose chemotherapy. *Pediatr. Blood Cancer* **63**, 1527–1534.
- Lex, A., Streit, M., Schulz, H.-J., Partl, C., Schmalstieg, D., Park, P.J., and Gehlenborg, N. (2012). *StratomeX*: visual analysis of large-scale heterogeneous genomics data for cancer subtype characterization. *Computer Graphics Forum* **31**, 1175–1184.
- Louis, D.N., Perry, A., Reifemberger, G., von Deimling, A., Figarella-Branger, D., Cavenee, W.K., Ohgaki, H., Wiestler, O.D., Kleihues, P., and Ellison, D.W. (2016). The 2016 World Health Organization classification of tumors of the central nervous system: a summary. *Acta Neuropathol.* **131**, 803–820.
- Maksimovic, J., Gordon, L., and Oshlack, A. (2012). *SWAN*: subset-quantile within array normalization for Illumina Infinium HumanMethylation450 BeadChips. *Genome Biol.* **13**, R44.
- Merico, D., Isserlin, R., Stueker, O., Emili, A., and Bader, G.D. (2010). Enrichment map: a network-based method for gene-set enrichment visualization and interpretation. *PLoS One* **5**, e13984.
- Mermel, C.H., Schumacher, S.E., Hill, B., Meyerson, M.L., Beroukhi, R., and Getz, G. (2011). *GISTIC2.0* facilitates sensitive and confident localization of the targets of focal somatic copy-number alteration in human cancers. *Genome Biol.* **12**, R41.
- Mo, Q., Wang, S., Seshan, V.E., Olshen, A.B., Schultz, N., Sander, C., Powers, R.S., Ladanyi, M., and Shen, R. (2013). Pattern discovery and cancer gene identification in integrated cancer genomic data. *Proc. Natl. Acad. Sci. USA* **110**, 4245–4250.
- Morrissy, A.S., Garzia, L., Shih, D.J., Zuyderduyn, S., Huang, X., Skowron, P., Remke, M., Cavalli, F.M., Ramaswamy, V., Lindsay, P.E., et al. (2016). Divergent clonal selection dominates medulloblastoma at recurrence. *Nature* **529**, 351–357.
- Northcott, P.A., Hielscher, T., Dubuc, A., Mack, S., Shih, D., Remke, M., Al-Halabi, H., Albrecht, S., Jabado, N., Eberhart, C.G., et al. (2011). Pediatric and adult sonic hedgehog medulloblastomas are clinically and molecularly distinct. *Acta Neuropathol.* **122**, 231–240.
- Northcott, P.A., Jones, D.T., Kool, M., Robinson, G.W., Gilbertson, R.J., Cho, Y.J., Pomeroy, S.L., Korshunov, A., Lichter, P., Taylor, M.D., and Pfister, S.M. (2012a). Medulloblastomics: the end of the beginning. *Nat. Rev. Cancer* **12**, 818–834.
- Northcott, P.A., Shih, D.J., Peacock, J., Garzia, L., Morrissy, A.S., Zichner, T., Stutz, A.M., Korshunov, A., Reimand, J., Schumacher, S.E., et al. (2012b). Subgroup-specific structural variation across 1,000 medulloblastoma genomes. *Nature* **488**, 49–56.
- Northcott, P.A., Lee, C., Zichner, T., Stutz, A.M., Erkek, S., Kawauchi, D., Shih, D.J., Hovestadt, V., Zapatka, M., Sturm, D., et al. (2014). Enhancer hijacking activates *GFI1* family oncogenes in medulloblastoma. *Nature* **511**, 428–434.
- Pei, Y., Liu, K.W., Wang, J., Garancher, A., Tao, R., Esparza, L.A., Maier, D.L., Udaka, Y.T., Murad, N., Morrissy, S., et al. (2016). HDAC and PI3K antagonists cooperate to inhibit growth of MYC-driven medulloblastoma. *Cancer Cell* **29**, 311–323.
- Ramaswamy, V., Northcott, P.A., and Taylor, M.D. (2011). FISH and chips: the recipe for improved prognostication and outcomes for children with medulloblastoma. *Cancer Genet.* **204**, 577–588.
- Ramaswamy, V., Remke, M., Bouffet, E., Faria, C.C., Perreault, S., Cho, Y.J., Shih, D.J., Luu, B., Dubuc, A.M., Northcott, P.A., et al. (2013). Recurrence patterns across medulloblastoma subgroups: an integrated clinical and molecular analysis. *Lancet Oncol.* **14**, 1200–1207.
- Ramaswamy, V., Remke, M., Bouffet, E., Bailey, S., Clifford, S.C., Doz, F., Kool, M., Dufour, C., Vassal, G., Milde, T., et al. (2016a). Risk stratification of childhood medulloblastoma in the molecular era: the current consensus. *Acta Neuropathol.* **131**, 821–831.
- Ramaswamy, V., Remke, M., Adamski, J., Bartels, U., Tabori, U., Wang, X., Huang, A., Hawkins, C., Mabbott, D., Laperriere, N., et al. (2016b). Medulloblastoma subgroup-specific outcomes in irradiated children: who are the true high-risk patients? *Neuro Oncol.* **18**, 291–297.
- Reimand, J., Arak, T., Adler, P., Kolberg, L., Reisberg, S., Peterson, H., and Vilo, J. (2016). *g:Profiler* — a web server for functional interpretation of gene lists (2016 update). *Nucleic Acids Res.* **44**, W83–W89.
- Remke, M., Hielscher, T., Northcott, P.A., Witt, H., Ryzhova, M., Wittmann, A., Benner, A., von Deimling, A., Scheurlen, W., Perry, A., et al. (2011). Adult medulloblastoma comprises three major molecular variants. *J. Clin. Oncol.* **29**, 2717–2723.
- Remke, M., Ramaswamy, V., Peacock, J., Shih, D.J., Koelsche, C., Northcott, P.A., Hill, N., Cavalli, F.M., Kool, M., Wang, X., et al. (2013). TERT promoter mutations are highly recurrent in SHH subgroup medulloblastoma. *Acta Neuropathol.* **126**, 917–929.
- Rutkowski, S., Cohen, B., Finlay, J., Luksch, R., Ridola, V., Valteau-Couanet, D., Hara, J., Garre, M.-L., and Grill, J. (2010). Medulloblastoma in young children. *Pediatr. Blood Cancer* **54**, 635–637.
- Shannon, P., Markiel, A., Ozier, O., Baliga, N.S., Wang, J.T., Ramage, D., Amin, N., Schwikowski, B., and Ideker, T. (2003). *Cytoscape*: a software environment for integrated models of biomolecular interaction networks. *Genome Res.* **13**, 2498–2504.
- Shen, R., Olshen, A.B., and Ladanyi, M. (2009). Integrative clustering of multiple genomic data types using a joint latent variable model with application to breast and lung cancer subtype analysis. *Bioinformatics* **25**, 2906–2912.
- Streit, M., Lex, A., Gratzl, S., Partl, C., Schmalstieg, D., Pfister, H., Park, P.J., and Gehlenborg, N. (2014). Guided visual exploration of genomic stratifications in cancer. *Nat. Methods* **11**, 884–885.
- Sturm, D., Witt, H., Hovestadt, V., Khuong-Quang, D.A., Jones, D.T., Konermann, C., Pfaff, E., Tonjes, M., Sill, M., Bender, S., et al. (2012). Hotspot mutations in H3F3A and IDH1 define distinct epigenetic and biological subgroups of glioblastoma. *Cancer Cell* **22**, 425–437.
- Taylor, M.D., Northcott, P.A., Korshunov, A., Remke, M., Cho, Y.J., Clifford, S.C., Eberhart, C.G., Parsons, D.W., Rutkowski, S., Gajjar, A., et al. (2012). Molecular subgroups of medulloblastoma: the current consensus. *Acta Neuropathol.* **123**, 465–472.

- Thompson, E.M., Hielscher, T., Bouffet, E., Remke, M., Luu, B., Gururangan, S., McLendon, R.E., Bigner, D.D., Lipp, E.S., Perreault, S., et al. (2016). Prognostic value of medulloblastoma extent of resection after accounting for molecular subgroup: a retrospective integrated clinical and molecular analysis. *Lancet Oncol.* *17*, 484–495.
- Wang, B., Mezlini, A.M., Demir, F., Fiume, M., Tu, Z., Brudno, M., Haibe-Kains, B., and Goldenberg, A. (2014). Similarity network fusion for aggregating data types on a genomic scale. *Nat. Methods* *11*, 333–337.
- Wilkerson, D.M., and Hayes, D.N. (2010). ConsensusClusterPlus: a class discovery tool with confidence assessments and item tracking. *Bioinformatics* *26*, 1572–1573.
- Zhao, F., Ohgaki, H., Xu, L., Giangaspero, F., Li, C., Li, P., Yang, Z., Wang, B., Wang, X., Wang, Z., et al. (2016). Molecular subgroups of adult medulloblastoma: a long-term single-institution study. *Neuro Oncol.* *18*, 982–990.
- Zhou, W., Laird, P.W., and Shen, H. (2016). Comprehensive characterization, annotation and innovative use of Infinium DNA Methylation BeadChip probes. *Nucleic Acids Res.* *45*, e22.
- Zhukova, N., Ramaswamy, V., Remke, M., Pfaff, E., Shih, D.J., Martin, D.C., Castelo-Branco, P., Baskin, B., Ray, P.N., Bouffet, E., et al. (2013). Subgroup-specific prognostic implications of TP53 mutation in medulloblastoma. *J. Clin. Oncol.* *31*, 2927–2935.
- Zou, H., and Hastie, T. (2005). Regularization and variable selection via the elastic net. *J. R. Stat. Soc. Ser. B* *67*, 301–320.

STAR★METHODS

KEY RESOURCES TABLE

REAGENT or RESOURCE	SOURCE	IDENTIFIER
Biological Samples		
763 primary medulloblastoma samples	This paper	N/A
Deposited Data		
Expression and methylation array raw and analyzed data	This paper	GEO: GSE85218
Expression array data (285 samples) (included as well in GSE85218)	Northcott et al., 2012b	GEO: GSE37382
SNP6 data	Northcott et al., 2012b	GEO: GSE37384
Oligonucleotides		
Primer for P53 see Table S5	Zhukova et al., 2013	N/A
TERT forward primer, 5'-CAG CGC TGC CTG AAA CTC-3'	Remke et al., 2013	N/A
TERT reverse primer, 5'-GTC CTG CCC CTT CAC CTT C-3'	Remke et al., 2013	N/A
Software and Algorithms		
Affy R Bioconductor package	Gautier et al., 2004	http://bioconductor.org/packages/release/bioc/html/affy.html
custom chip definition file (CDF) hugene11sthsensgcdf (v19.0.0).	Dai et al., 2005	http://brainarray.mbni.med.umich.edu/Brainarray/Database/CustomCDF/19.0.0/ensg.asp
arrayQualityMetrics R Bioconductor package (v3.22.0)	Kauffmann et al., 2009	https://www.bioconductor.org/packages/release/bioc/html/arrayQualityMetrics.html
minfi R Bioconductor package (v1.6.0) including SWAN normalization method	Aryee et al., 2014; Maksimovic et al., 2012	http://bioconductor.org/packages/release/bioc/html/minfi.html
NMF R package (v0.20.6)	Gaujoux and Seoighe, 2010	https://cran.r-project.org/web/packages/NMF/index.html
conumee R Bioconductor package (v0.99.4)	Hovestadt et al., 2013; Sturm et al., 2012	http://bioconductor.org/packages/release/bioc/html/conumee.html
GISTIC2 method (v6.2)	Mermel et al., 2011	http://portals.broadinstitute.org/cgi-bin/cancer/publications/pub_paper.cgi?mode=view&paper_id=216&p=t
ConsensusClusterPlus R Bioconductor package (v1.24.0)	Wilkerson and Hayes, 2010	https://www.bioconductor.org/packages/release/bioc/html/ConsensusClusterPlus.html
SNFtool R package (v2.2.0)	Wang et al., 2014	https://cran.r-project.org/web/packages/SNFtool/index.html
MethylMix R Bioconductor package (2.0.0)	Gevaert, 2015	https://www.bioconductor.org/packages/release/bioc/html/MethylMix.html
Infinium DNA Methylation BeadChip (450K) probe annotation on hg38	Zhou et al., 2016	http://zwdzwd.github.io/InfiniumAnnotation
StratomeX tool as part of the Caleydo suite (v3.1.5)	Streit et al., 2014; Lex et al., 2012	http://caleydo.org/tools/stratomeX/
g:profiler	Reimand et al., 2016	http://biit.cs.ut.ee/gprofiler/
Cytoscape (v3.2.0)	Shannon et al., 2003	http://www.cytoscape.org/
Cytoscape Enrichment map	Merico et al., 2010	http://apps.cytoscape.org/apps/enrichmentmap
IClusterPlus R Bioconductor package (v1.10.0)	Mo et al., 2013	https://www.bioconductor.org/packages/release/bioc/html/iClusterPlus.html

CONTACT FOR REAGENT AND RESOURCE SHARING

Further information and requests for resources and reagents should be directed to and will be fulfilled by the Lead Contact, Michael D Taylor (mdtaylor@sickkids.ca).

EXPERIMENTAL MODEL AND SUBJECT DETAILS

Acquisition of Patient Samples

All medulloblastoma samples were collected at diagnosis after obtaining informed consent from subjects as part of the Medulloblastoma Advanced Genomics International Consortium. Approval was obtained from institutional research ethics boards at the following institutions: The Hospital for Sick Children, Children's Hospital of Pittsburgh, Seoul National University Children's Hospital, The Children's Memorial Health Institute, Institute of Pediatric Hematology and Oncology, Mayo Clinic, The Chinese University of Hong Kong, John Hopkins University School of Medicine, University of Alabama at Birmingham, Seattle Children's Hospital, University of California San Francisco, Burdenko Neurosurgical Institute, McMaster University, Erasmus University Medical Center, Asan Medical Center, Kitasato University School of Medicine, Hospital Pediatría Centro Médico Nacional Century XXI, Masaryk University, Fondazione IRCCS Istituto Nazionale Tumori, Emory University, Osaka National Hospital, University of Debrecen, University of Naples, Washington University School of Medicine, Montreal Children's Hospital, Hospital Sant Joan de Déu, Virginia Commonwealth University, Chonnam National University Hwasun Hospital and Medical School, Children's Health Queensland Hospital and Health Service, University of Calgary, University of Sao Paulo, Cincinnati Children's Hospital Medical Center, Hospital de Santa Maria, Lisbon, University of Arkansas for Medical Sciences, Catholic University Medical School, David Geffen School of Medicine at UCLA, The University of Sydney, Kumamoto University Graduate School of Medical Science, Saint Louis University School of Medicine, Hospital Infantil de Mexico Federico Gomez, Rainbow Babies & Children's Hospital. Patients were selected only if their treatment plan required surgical resection. Samples were obtained as fresh frozen tissue from the time of diagnosis and stored at -80°C until processed for the purification of nucleic acids. Tumor isolates were partitioned for both DNA and RNA extraction. Using all information in our hands, we selected only primary tumor medulloblastoma samples for this study and removed duplicates. The sex and gender of the 763 medulloblastoma patients used in this study are presented in [Table S1](#).

METHOD DETAILS

Nucleic Acid Extraction

DNA extraction was performed by incubation with proteinase K overnight at 55°C followed by three sequential phenol extractions and ethanol precipitation. Total RNA was isolated using the TriZol method where tissue was homogenized in a Precellys 24 tissue homogenizer (Bertin Technologies, France) in Trizol using strict RNAase free conditions. DNA was quantified using the Picogreen method and RNA quantified using a NanoDrop 1000 instrument (Thermo Scientific) and integrity assessed by agarose gel electrophoresis (DNA) or Agilent 2100 Bioanalyzer (RNA) at the Centre for Applied Genomics (TCAG) at the Hospital for Sick Children in Toronto, Canada. RNA with an RNA integrity number of 7 or higher was required for analysis by Affymetrix Gene Arrays.

Expression and Methylation Data

To generate gene expression array profiling, 400ng of total RNA was processed and hybridized to the Affymetrix Gene 1.1 ST array at the Centre for Applied Genomics (TCAG) at the Hospital for Sick Children (Toronto, Canada) according to manufacturers instructions. In addition, all samples were analyzed on the Illumina Infinium HumanMethylation450 BeadChips at TCAG (Toronto, ON) according to manufacturer's instructions.

TERT Promoter and TP53 Sequencing

TERT promoter mutational status was determined using direct sanger sequencing and genotyping as previously described where sufficient DNA was available ([Remke et al., 2013](#)). Two primers (forward primer, 5'-CAG CGC TGC CTG AAA CTC-3'; reverse primer, 5'-GTC CTG CCC CTT CAC CTT C-3') were designed to amplify a 163-bp product encompassing C228T and C250T hotspot mutations in the *TERT* promoter—corresponding to the positions 124 and 146 bp, respectively, upstream of the ATG start site. Two fluorescent LNA probes were designed with different fluorescent dyes to allow single-tube genotyping. One probe was targeted to the WT sequence (*TERT* WT, 5'-5HEX-CCC CTC CCG G-3IABkFQ-3'), and one was targeted to either of the two mutations (*TERT* mut, 5'-56FAM-CCC CTT CCG G-3IABkFQ). Primer and probes were custom designed by Integrated DNA Technologies (Coralville, Iowa, USA) using internal SNP design software, and sequence homogeneity was confirmed by comparison to all available sequences on the GenBank database using BLAST (). Primers were optimized to avoid for hairpins and homo- and heterodimers. Primers and probes were obtained from Integrated DNA Technologies.

Real-time PCR was performed in 25 μl reaction mixtures containing 12.5 μl of TaqMan Universal Master Mix II with UNG (Applied Biosystems), 900 nM concentrations of each primer, 250 nM *TERT* WT probe, 250 nM *TERT* MUT probe, and 1 μl (25 ng) of sample DNA. Thermocycling was performed on the StepOnePlus (Applied Biosystems) and consisted of 2 min at 50°C , 10 min at 95°C , and 40 cycles of 95°C for 15 s and 60°C for 1 min.

Analysis was performed using StepOne Software, version 2.1. Samples were considered mutant if they had CT values of ≤ 39 cycles. Each sample was verified visually by examining the PCR curves generated to eliminate false positives due to aberrant light emission. End-point allelic discrimination genotyping was performed by visually inspecting a plot of the fluorescence from the WT probe versus the MUT probe generated from the post-PCR fluorescence read.

TP53 mutational status was determined using direct sanger sequencing as previously described where sufficient DNA was available (Zhukova et al., 2013). We used amplitaq gold and after purification with ampure beads, forward and reverse sequencing primers using dGTP BigDye Terminator v3.0 Cycle Sequencing Ready Reaction Kit (Life Technologies), and 5 % DMSO on the ABI3730XL capillary genetic analyzer (Life Technologies). The sequencing primers are the same as the PCR primers. The *TP53* primers along with the melting temperature are presented in Table S5.

QUANTIFICATION AND STATISTICAL ANALYSIS

Microarray Gene Expression Analysis

To generate gene expression array profiling, 400ng of total RNA was processed and hybridized to the Affymetrix Gene 1.1 ST array at TCAG according to manufacturer's instructions. Two hundred and eighty-five arrays were previously generated (GEO accession GSE37382) and included in the analysis. Expression data were analyzed in the R environment (v3.1.1). We used the affy package (v1.44.0) (Gautier et al., 2004) and the custom chip definition file (CDF) hugene11sthsensgcdf (v19.0.0).

(<http://brainarray.mbni.med.umich.edu/Brainarray/Database/CustomCDF/19.0.0/ensg.asp>) (Dai et al., 2005) to load and summarize the expression of 21,641 Ensembl (ENSG) genes and process the data. Samples flagged by the arrayQualityMetrics Bioconductor package (v3.22.0) (Kauffmann et al., 2009) were removed due to low quality. Expression data were normalized using the rma method.

Unsupervised clustering using NMF using top 10,000 most variably expressed genes (determined by the standard deviation) was carried out using the NMF package (v0.20.6) (Gaujoux and Seoighe, 2010). We reselected the top most 10,000 variably expressed genes for each subset of samples on which we ran NMF.

Genome Wide Methylation Analysis

All samples were analyzed on the Illumina Infinium HumanMethylation450 BeadChips at TCAG (Toronto, Ontario) according to manufacturer's instructions. Bisulfite conversion was performed using the EZ DNA Methylation™ Kit (Zymo, Irvine, CA). Samples were processed as per manufacturer's instructions. Raw data files (.idat) generated by the Illumina iScan array scanner were processed in the R statistical environment (v3.0.0 and 3.1.1) using the minfi (v1.6.0) (Aryee et al., 2014) and IlluminaHumanMethylation450kmanifest (v0.4.0) R Bioconductor packages. We checked all samples for unexpected genotype matches by pairwise correlation of the 65 genotyping probes on the 450k arrays, allowing us to remove remaining duplicates. We ran the detectionP function from the minfi package to identify probes and samples with low quality. Samples were removed if more than 1% of their probes had a p value above 0.01 and probes were removed if their p value was above 0.01 in at least 5% of samples. We removed probes on sex chromosomes as well as those located on or close to known single nucleotide polymorphisms (SNP). We retained a total of 321,174 probes for the analysis. The data was normalized using the SWAN method as part of the minfi package (Maksimovic et al., 2012). We generated both the beta and logitB values matrix values. Unsupervised clustering using the top 10,000 most variably methylated probes defined by the standard deviation was carried out using the NMF package (v0.20.6). We reselected the top most 10,000 variably methylated probes for each subset of samples on which we ran NMF.

Methylation Array Copy Number Analysis

Copy number inference from methylation arrays and identification of recurrent broad events. Copy number segmentation was performed from genome wide methylation arrays using the conumee package (v0.99.4) in the R statistical environment (v3.2.3) as previously described (Hovestadt et al., 2013; Sturm et al., 2012). Segment files were generated for each subgroup and subtype.

Identification of recurrent broad copy number events (arm level chromosomal events) was performed from segmented copy number derived from methylation data (as described above). The log₂ R ratio (LRR) of each chromosome was calculated using a size-weighted mean of all segments mapping to the chromosome. A chromosome was declared gained if its LRR was greater than 0.2, lost if the LRR was less than -0.2, and balanced otherwise. Unlike GISTIC, gained and lost broad events were analyzed together. The significance of the frequency of each broad event was tested using the exact binomial test. Each broad event frequency was compared to the background frequency, which was determined from a robust regression of the observed frequencies with respect to gene content (i.e. number of RefSeq genes) across all chromosomes. This approach was motivated by GISTIC's broad event analysis.

SNP6 Copy Number Analysis

Affymetrix SNP6 CEL files were processed as previously described (Northcott et al., 2012b) (GEO accession GSE37384). Copy number states were estimated as described previously using the hg18 reference genome. Segmented copy number estimates from SNP6 arrays were processed for input with the GISTIC2 method (v6.2) using the default parameters (Mermel et al., 2011) for the identification of recurrent focal copy number events.

Clinical Correlation and Survival Analysis

Progression-free survival and overall survival was right-censored at 5 years and analyzed by the Kaplan-Meier method and p value were reported using the log-rank test. Associations between covariates and risk groups were tested by the Fisher's exact test. Continuous variables were tested using non-parametric measures, specifically the Mann-Whitney U test or Kruskal-Wallis test. The significance of chromosome arm frequencies were evaluated using the exact binomial test, comparing the observed frequency to the expected frequency derived from a robust regression of event frequency and gene content, in a similar manner to the 'broad analysis' in GISTIC2. All statistical analyses were performed in the R statistical environment (v3.2.3), using R packages survival (v2.37-7), and ggplot2 (v1.0.0).

Group 3 and Group 4 Analysis

K-means clustering was performed using the top 10,000 most variable methylation probes (determined by median absolute deviation) of Group 3 and Group 4 samples (n=470). Consensus clustering was obtained using the ConsensusClusterPlus R Bioconductor package (v1.24.0) (Wilkerson and Hayes, 2010) with 1,000 repetitions in the R statistical environment (v3.2.2). Similar approach was used on the top 10,000 most variable genes of this set of 470 Group 3 and Group 4 samples. In addition, the NMF method was run (as described above) for both expression and methylation data on the same set of Group 3 and Group 4 samples.

We identified the outlier samples moving from Group 3 to Group 4 from the gene expression and DNA methylation NMF results using the following rule. We identified at k=2 the Group 3 and Group 4 clusters using the known subgroups of the samples (each group had a larger proportion of samples of a particular subgroup). At k=3, we identified which cluster(s) are largely composed of Group 3 and Group 4 (two Group 3 and one Group 4 clusters for the expression data, and one Group 3 and two Group 4 clusters for the DNA methylation data, Figure S1G). Then we counted the number of samples that were initially considered to be of a particular subgroup for k=2 and moved to be in another subgroup at k=3 (Figure 1D). Similar approach has been used to detect the outlier samples moving from Group 3 to Group 4 in the k-means consensus clustering (Figures 1D and S1H).

Similarity Network Fusion Analysis (SNF)

The Similar Network Fusion (SNF) method was run on 763 primary tumor samples using both gene expression and DNA methylation data (Wang et al., 2014). The SNF method does not require any prior feature selection so we used the full matrix of gene expression (21,641 genes) and the full matrix of methylation data (logitB values, 321,174 probes). We used the SNFtool R package (v2.2.0) with the parameters $K = 50$, $\alpha = 0.6$, $T = 50$. Spectral clustering implemented in the SNFtool package was run on the SNF fused similarity matrix to obtain the groups corresponding to k=2 to 20.

We obtained four cluster at k=4 corresponding to the four medulloblastoma subgroups; WNT (n= 70), SHH (n=233), Group 3 (n=144), Group 4 (n=326). For each of these four subgroups we then ran the SNF method independently with the following parameters and clustered the resulting fused similarity matrix with spectral clustering using k=2 to 8.

Parameters:

WNT: $K = 10$, $\alpha = 0.6$, $T = 50$

SHH: $K = 40$, $\alpha = 0.6$, $T = 50$

Group 3: $K = 40$, $\alpha = 0.6$, $T = 50$

Group 4: $K = 60$, $\alpha = 0.6$, $T = 50$

Group 3 and Group 4: $K = 80$, $\alpha = 0.6$, $T = 50$

We identified the top associated genes and methylation probes that have the largest agreement with the final fused network structure. To do so we computed the Normalized Mutual Information (NMI) score (as part of the SNFtool package) for each feature (i.e each gene and methylation probe). For each feature, we constructed a patient network based on the feature alone and subsequently used spectral clustering. We then compared the result of the resultant clustering to the one obtained from the whole fused similarity matrix by computing the NMI score as previously described (Wang et al., 2014). As mentioned in this paper, a score of 1 indicates the strongest feature and shows that the network of patients based on the given feature leads to the same groups as the fused network. A score of 0 means that there is no agreement between the groups that can be derived from the feature and the fused network groups. We therefore ranked all features according to their NMI scores that represent their importance for the fused network. We then selected a list of top 1% and top 10% features (also called associated genes and methylation probes) for each dataset (Figure 2 and S2A–S2D) for subsequent analysis. Those top features have expression or methylation patterns that are the most informative when determining our final subtypes using individual features.

Groups Visualization Using StratomeX

We used the StratomeX tool as part of the Caleydo suite (v3.1.5) to visualize the grouping of samples and the relationship between the groups resulting from different datasets, methods and/or parameterization of clustering (Streit et al., 2014; Lex et al., 2012). Sample group labels obtained by spectral clustering of the SNF fused similarity matrix or independent NMF clustering was imported to StratomeX software. The groups were colored according to the subgroup or subtype (if any) and reordered to show the relationship between the different clustering results (columns). In this study, we only used StatromeX for visualization and did not use its analytical functionalities.

Network Visualization with Cytoscape

From the fused similarity matrix returned by the SNF method, we retrieved all the patient pairs for which the values (W) was superior to the median values of all W pairs and imported those paired in Cytoscape (v3.2.0) (Shannon et al., 2003). We used the edged-weighted Spring embedded layout with the W values for visualization showing the edges in Figure 1B, and hiding the edges to only show the nodes (i.e patients) for Figures 3A, 4A, 5A, and 6A.

Relationship between Associated Genes and Probes

We evaluated the relationship between the gene expression features and the DNA methylation probe features in each subgroup. We applied the MethylMix R Bioconductor package (Gevaert, 2015) developed to identify potential cancer driver genes affected by hypo or hypermethylation changes, i.e. looking for anti-correlation between the methylation level and gene expression levels across samples. We obtained the probes annotations for hg38 from Zhou et al. (Zhou et al., 2016, Online supplemental data, <http://zwdzwd.github.io/InfiniumAnnotation>). We focused on probes within 1500 bp of the transcription start site (TSS) and identify 1342, 1573, 1673 and 1673 candidate driver genes for WNT, SHH, Group 3 and Group 4, respectively. Among those, 8, 18, 13 and 28 WNT, SHH, Group 3 and Group 4 genes, respectively, where in our features genes and had anti-correlated probes present in the top DNA methylation features, representing therefore only 3.7, 8.3, 6 and 13% of the feature genes (Figure 2C).

Pathway Enrichment Analysis

Pathway enrichment analysis was performed with g:Profiler and visualized as Enrichment Map in Cytoscape (Reimand et al., 2016; Merico et al., 2010; Shannon et al., 2003). We considered the top 10% associated genes (as described above) that were the most relevant for the final subtypes. For each subtype, we ranked up-regulated genes by their z-scores and analyzed the resulting gene lists with the ordered query setting of g:Profiler using pathways and processes with more than 5 and up to 1000 genes. Multiple testing correction was conducted with the default method of g:Profiler. Biological processes from the Gene Ontology, pathways from Reactome and KEGG, and protein complexes from CORUM were included in the enrichment analysis and other data sources were excluded. Electronic annotations (IEA) from Gene Ontology were excluded to only cover high-confidence gene annotations. Processes and pathways with g:profiler FDR corrected q values <0.05 were considered significant. Enriched categories were further filtered: pathways and processes with less than three associated genes were discarded.

Enrichment maps represent biological processes and pathways enriched in subtype-specific up-regulated genes. Each node represents a process or pathway; nodes with many shared genes are grouped and labeled by biological theme. Nodes sizes are proportional to the number of genes in each process, in each subgroup. Process and pathways connected by edges have genes in common, shorter edges represent stronger edges with Jaccard and Overlap coefficient combined by the Enrichment Map app of Cytoscape at cutoff value 0.66. Nodes are colored according to the subtype in which the process is enriched; processes enriched in more than one subtype have multiple colors.

Enrichment map visualization was manually curated to group functionally similar groups of pathways and to remove redundant groups and singletons. Connected nodes and unconnected but actionable nodes are shown.

Classifier Description

In this study, we used seven classifiers based on diverse machine learning approaches. Ridge logistic regression (labeled as *Ridge LR*) is a regression model, assigning weight to each feature to make a binary prediction. L2 regularization shrinks the weights to avoid overfitting. Lasso logistic regression (labeled as *Lasso LR*) works the same way as Ridge but uses L1 regularization instead, which sets some of the weights to zero, effectively performing feature selection again to avoid overfitting. Elastic net logistic regression (labeled as *Elastic Net*) also works similarly to Ridge but uses a linear combination of L1 and L2 regularizations and is able to select correlated features (through L2) while still performing feature selection (setting some of the weights to zero) through L1. Decision tree (labeled as *Decis. Tree*) utilizes a tree-structured graph with inner nodes representing decision rules and end nodes representing the classification decisions. Each path from root to a leaf in such a tree represents a classification rule. The individual decision rules are selected according to the information gain criterion. Random forest (labeled as *Rand. Forest*) uses an ensemble of decision trees to make classification predictions. Each decision tree uses a random subset of features trained on a bootstrapped set of samples. The output is the mode of the classification from all decision trees in the random forest. Support Vector Machines (SVM) make classification predictions by first transforming the data according to a chosen kernel and then constructing a maximum margin classifier such that the different classes are separated by the decision hyperplane as much as possible. SVM with linear kernel (labeled as *SVM lin*) performs a linear transformation of the data, whereas SVM with radial basis function kernel (labeled as *SVM rbf*) performs a Gaussian transformation of the data.

Prior to training any of the classifiers, we used Kruskal-Wallis H test, also know as “one-way ANOVA on ranks” to constrain the feature space. This way we selected top 1% of genes (resulting in 216 genes) and top 1% of CpG methylation probes (resulting in 3212 methylation probes) whose expression and methylation, respectively, is most predictive of the cluster assignment (done by the spectral clustering on the SNF fused similarity matrix) of samples in the training set. This feature selection procedure was repeated in each training / testing split of the data set, using the training set only.

All analyses were performed in R version 3.2.3. We used the `glmnet` package for Elastic Net, Lasso, and Ridge; the `rpart` package for Decision tree; the `randomForest` package for Random forest, the `kernLab` package for SVM. Training of Random forest and SVM was done using the `caret` package. The AUPRC (area under the precision-recall curve) values were calculated using the `PRROC` package.

Training and Selection of the Classifiers

We performed five classification tasks: the medulloblastoma subgroup classification, and then subtype classification within each of the four medulloblastoma subgroups. Cluster assignment by spectral clustering on the SNF fused similarity matrix was taken as the ground truth label assignment for the study cohort subgroup and subtype classification. For each of these tasks we trained 7 classification models using the concatenation of the top 1% expression and the top 1% methylation features as feature set.

For each task, we split the study cohort data set randomly to 70% training set and 30% testing set splits. The top 1% feature selection procedure, as described above, was then run on the training set. The selected features were used for both training and testing of the classifiers. Individual classification models were subsequently trained in 5-fold cross validation on the training set. On the testing set we measured the performance of the classifiers in the terms of classification accuracy and the area under the precision-recall curve (AUPRC). The entire described procedure was repeated 100 times. We report AUPRC and accuracy means and standard deviation over these 100 runs, as well as the average percentages of subtype predicted for the reference subtypes (Table S2).

COCA Analysis

We performed the COCA analysis as described in the TCGA pan cancer paper (Hoadley et al., 2014). We applied NMF clustering on gene expression and DNA methylation data individually for each subgroup. We proceeded to create a matrix of 0 and 1 with all the samples (as column) and the different groups (as row, one row per group obtained for each clustering). 1 indicated the presence of a sample in a group. This matrix was then clustered with k-means consensus clustering as performed in the TCGA pan cancer paper.

iCluster Analysis

We also performed multi-platform clustering using iCluster. We applied the R Bioconductor iClusterPlus package (the newest version of iCluster) to perform the analysis (Shen et al., 2009; Mo et al., 2013). It is necessary to select a set of features for each dataset to run iCluster. We tested selection of features on the maximum variance and on the MAD (median absolute deviation) and different percentages of features. We selected the top 15% most variable expressed genes and 1% most variable methylated probes as defined by median absolute deviation. We chose these numbers to allow for an equal representation of variable genes and methylated probes, which resulted in 3000 features per dataset. We performed multiple clusterings with different values of the lambda parameter and settled on $\lambda = 1$. After performing the clustering, we confirm that almost all features are used to some degree in the model which implies that the results are not entirely driven by one data type. Indeed, this reassured us that the parameters we selected would allow for a robust multi-platform integrated analysis.

When applying iCluster across the entire dataset, we are unable to recover the 4 subgroups of medulloblastoma at $k=4$. When comparing the demographics of these 4 groups, and cross referencing to the SNF subgroups, WNT and Group 3 cluster together with two SHH groups emerging. When we increased to 5 groups we are able to clearly split WNT and Group 3. To determine the congruence between iCluster and SNF in defining the subtypes of subgroups, we first defined the subtypes using SNF, and then applied iCluster individually to each subgroup. Overall the subtypes as defined by iCluster were in agreement with SNF/Spectral clustering groups in 72-84% of instances. When we take into account that in some instances, iCluster recovered similar subtypes at a different number of groups, then the agreement increases to 74-90% (for example, some groups in iCluster split in two but correspond strongly to one cluster by SNF).

DATA AND SOFTWARE AVAILABILITY

The expression and methylation array data has been deposited in GEO under the accession number GSE85218. The previously published data is available in GEO under the accession numbers GSE37382 and GSE37384.