

## SUPPLEMENTARY INFORMATION

### Supplementary Methods

**Array Data Analysis.** Following hybridization, each chip passed quality assurance and control procedures using the Affymetrix quality control algorithms provided in Expressionist Refiner module (Genedata AG, Basel, Switzerland)<sup>4</sup>. Probe signal levels were quantile-normalized and summarized using the GeneChip - Robust Multichip Averaging (CG-RMA) algorithm<sup>1</sup>. Normalized files were imported into Partek (Partek Inc., St. Louis, MO, USA) and into the Expressionist Analyst module for principal component analysis (PCA), unsupervised clustering, and to assess significant differences in gene expression<sup>5,6</sup>. The data generated by both programs showed consistent patterns. To develop sample size estimates for gene expression profiling, we started with theoretical principles and then applied empirical observations to support the sample size for these experiments *a priori*. The Canine\_2.0 gene expression chip contains ~43,000 annotated sequences derived from the 7.56x canine genome<sup>2</sup>. These represent virtually every known gene and a complement of expressed sequence tags that provide strong redundancy for expression profiling. We next considered that False Discovery Rate statistical analysis provided a suitable method to set thresholds for significance of elevated or reduced gene expression, but additional multivariate analyses and gene set enrichment would add further value to the analysis. We anticipated the data might not be normally distributed; so, non-parametric tests might be needed. As there is no analytical estimate of the power of the Kruskal-Wallis test after false discovery rate corrections, an approximation is useful in the case of small sample sizes. We can estimate the proportion of times when perfect rank separation between conditions might occur by chance as

$2N!N!/(2N)!$ , where  $N$  is the number of samples in each group. The Power Atlas (<http://www.poweratlas.org/>), allowed us to obtain an empirical estimate that the imbalanced sample sets used for these experiments should provide >90% power at  $p = 0.05$  to identify true positives, although the power to identify true negatives could be lower.

The correlation coefficient ( $r^2$ ) for expression values of all probes between the duplicated samples was >0.95. Probe IDs were mapped to corresponding canine Entrez Gene IDs using Affymetrix NetAffx EntrezGene Annotation. Prior to hierarchical clustering, normalized chip data were median-centered and  $\log_2$ -transformed.

Unsupervised clustering was done using Gene Cluster 3.0 for Mac OS X (C Clustering Library 1.47) with correlation based on average linkage. Gene Cluster 3.0 data were visualized in Java TreeView<sup>3</sup>. Two group t-tests were done to determine genes that were differentially expressed between groups. As with all microarray analysis, correction for multiple testing is required. We further selected for driver genes with the largest effect by restricting analysis to differentially expressed genes that showed large fold changes (>3) and highly significant  $p$ -values <0.001. Though batch effects were discernible between the two cohorts, they did not affect analysis as each cohort was analyzed independently, with cohort-1 used as a training-set and cohort-2 used as a validation-set. Gene expression data were deposited in Gene Expression Omnibus (GEO).

### **Network identification and canonical pathway analysis of differentially expressed**

**genes.** Ingenuity Pathway Analysis (IPA) software (Ingenuity Systems, Redwood City, CA) was used to define functions and canonical pathways specifically enriched in the sets of genes using BH multiple testing corrections to assess significance <sup>4</sup>. Gene Set Enrichment Analysis (GSEA, <http://www.broad.mit.edu/gsea/>) was similarly used to define enriched functional pathways as described previously <sup>6</sup>. Statistical significance was estimated using phenotype-based permutations, with the attained *p*-values adjusted for multiple hypothesis testing.

### Supplementary References

- 1 Irizarry RA, Wu Z, Jaffee HA: Comparison of Affymetrix GeneChip expression measures. *Bioinformatics*. 2006 22 (7): 789-794.
- 2 Lindblad-Toh K, Wade CM, Mikkelsen TS, et al.: Genome sequence, comparative analysis and haplotype structure of the domestic dog. *Nature*. 2005 438 (7069): 803-819.
- 3 Saldanha AJ: Java Treeview--extensible visualization of microarray data. *Bioinformatics*. 2004 20 (17): 3246-3248.
- 4 Scott MC, Sarver AL, Gavin KJ, et al.: Molecular subtypes of osteosarcoma identified by reducing tumor heterogeneity through an interspecies comparative approach. *Bone*. 2011 49 (3): 356-367.
- 5 Tamburini BA, Phang TL, Fosmire SP, et al.: Gene expression profiling identifies inflammation and angiogenesis as distinguishing features of canine hemangiosarcoma. *BMC Cancer*. 2010 10 (1): 619.
- 6 Tamburini BA, Trapp S, Phang TL, et al.: Gene expression profiles of sporadic canine hemangiosarcoma are uniquely associated with breed. *PLoS ONE*. 2009 4 (5): e5549.

Supplementary Tables

**Table S1a.** Functional pathways enriched in T-cell lymphoma (vs. B-cell lymphoma)

<b>Functions</b>	<b>p-Value</b>
Proliferation of lymphocytes	8.95E-25
Proliferation of immune cells	4.21E-24
Proliferation of T lymphocytes	3.41E-23
Immune response	1.79E-22
Cell death of immune cells	2.99E-22

**Table S1b.** Functional pathways enriched in B-cell lymphoma (vs. T-cell lymphoma)

<b>Functions</b>	<b>p-Value</b>
Activation of B lymphocytes	7.84E-17
Developmental process of B lymphocytes	1.06E-15
Quantity of B lymphocytes	1.51E-15
Proliferation of B lymphocytes	2.60E-15
Antibody response	5.35E-12

**Table S1c.** Functional pathways enriched in high-grade T-cell lymphoma (vs. low-grade T-cell lymphoma)

<b>Functions</b>	<b>p-Value</b>
Cell division process of chromosomes	1.58E-23
Segregation of chromosomes	2.29E-18
Mitosis	1.11E-13

Ploidy	3.02E-11
Cell cycle progression	5.18E-11

**Table S1d.** Functional pathways enriched in low-grade B-cell lymphoma (vs. high-grade B-cell lymphoma)

<b>Functions</b>	<b>p-Value</b>
Survival of T lymphocytes	3.19E-08
Activation of cells	8.58E-08
Survival of lymphocytes	2.52E-07
Survival of blood cells	2.54E-07
Cell death of immune cells	2.60E-07

### **Figure Legends**

**Figure s1.** a; Principal component analysis of normalized gene expression profiles of canine high- and low-grade T-cell and B-cell lymphoma reveals three molecular groups in canine lymphoma samples. Molecular relatedness of samples is described by distance in three-dimensional space as defined by three principal components of molecular variability. Each identifiable molecular group was labeled with a letter (a, b, c) for ease of identification. b; Heat map showing expression data for genes (N=859) with variance >1.0 and >8 fold change in at least 3 profiles. Colors represent median-centered fold change expression following  $\log_2$  transformation (a quantitative representation of the

colors is provided in the scale at the bottom). Upregulated genes are shown in red and downregulated genes are shown in green.

**Figure s2.** Immunophenotyping of canine B-cell lymphoma. Forty-eight independent samples of canine B-cell lymphoma (28 high-grade tumors and 20 low-grade tumors) were immunophenotyped by flow cytometry using antibodies against CD3, CD5, CD21, and CD22. The frequency of T cells and B cells was enumerated from analysis of >10,000 cells per sample. The box-plot provides a visual summary of the data. The two groups were statistically significantly different ( $p=0.0079$ ).