



# Genome-wide analysis of cancer/testis gene expression

## Citation

Hofmann, O., O. L. Caballero, B. J. Stevenson, Y.-T. Chen, T. Cohen, R. Chua, C. A. Maher, et al. 2008. Genome-Wide Analysis of Cancer/testis Gene Expression. *Proceedings of the National Academy of Sciences* 105, no. 51: 20422–20427. doi:10.1073/pnas.0810777105.

## Published Version

doi:10.1073/pnas.0810777105

## Permanent link

<http://nrs.harvard.edu/urn-3:HUL.InstRepos:29074749>

## Terms of Use

This article was downloaded from Harvard University's DASH repository, and is made available under the terms and conditions applicable to Other Posted Material, as set forth at <http://nrs.harvard.edu/urn-3:HUL.InstRepos:dash.current.terms-of-use#LAA>

## Share Your Story

The Harvard community has made this article openly available.  
Please share how this access benefits you. [Submit a story](#).

[Accessibility](#)

# Genome-wide analysis of cancer/testis gene expression

Oliver Hofmann<sup>a,b,1</sup>, Otavia L. Caballero<sup>c</sup>, Brian J. Stevenson<sup>d,e</sup>, Yao-Tseng Chen<sup>f</sup>, Tzeela Cohen<sup>c</sup>, Ramon Chua<sup>c</sup>, Christopher A. Maher<sup>b</sup>, Sumir Panji<sup>b</sup>, Ulf Schaefer<sup>b</sup>, Adele Kruger<sup>b</sup>, Minna Lehtvaslaiho<sup>b</sup>, Piero Carninci<sup>g,h</sup>, Yoshihide Hayashizaki<sup>g,h</sup>, C. Victor Jongeneel<sup>d,e</sup>, Andrew J. G. Simpson<sup>c</sup>, Lloyd J. Old<sup>c,1</sup>, and Winston Hide<sup>a,b</sup>

<sup>a</sup>Department of Biostatistics, Harvard School of Public Health, 655 Huntington Avenue, SPH2, 4th Floor, Boston, MA 02115; <sup>b</sup>South African National Bioinformatics Institute, University of the Western Cape, Private Bag X17, Bellville 7535, South Africa; <sup>c</sup>Ludwig Institute for Cancer Research, New York Branch at Memorial Sloan-Kettering Cancer Center, 1275 York Avenue, New York, NY 10021; <sup>d</sup>Ludwig Institute for Cancer Research, Lausanne Branch, 1015 Lausanne, Switzerland; <sup>e</sup>Swiss Institute of Bioinformatics, 1015 Lausanne, Switzerland; <sup>f</sup>Weill Medical College of Cornell University, 1300 York Avenue, New York, NY 10021; <sup>g</sup>Genome Exploration Research Group (Genome Network Project Core Group), RIKEN Genomic Sciences Center (GSC), RIKEN Yokohama Institute, 1-7-22 Suehiro-cho, Tsurumi-ku, Yokohama, Kanagawa, 230-0045, Japan; and <sup>h</sup>Genome Science Laboratory, Discovery Research Institute, RIKEN Wako Institute, 2-1 Hirosawa, Wako, Saitama, 3510198, Japan

Contributed by Lloyd J. Old, October 28, 2008 (sent for review June 6, 2008)

**Cancer/Testis (CT) genes, normally expressed in germ line cells but also activated in a wide range of cancer types, often encode antigens that are immunogenic in cancer patients, and present potential for use as biomarkers and targets for immunotherapy. Using multiple in silico gene expression analysis technologies, including twice the number of expressed sequence tags used in previous studies, we have performed a comprehensive genome-wide survey of expression for a set of 153 previously described CT genes in normal and cancer expression libraries. We find that although they are generally highly expressed in testis, these genes exhibit heterogeneous gene expression profiles, allowing their classification into testis-restricted (39), testis/brain-restricted (14), and a testis-selective (85) group of genes that show additional expression in somatic tissues. The chromosomal distribution of these genes confirmed the previously observed dominance of X chromosome location, with CT-X genes being significantly more testis-restricted than non-X CT. Applying this core classification in a genome-wide survey we identified >30 CT candidate genes; 3 of them, PEPP-2, OTOA, and AKAP4, were confirmed as testis-restricted or testis-selective using RT-PCR, with variable expression frequencies observed in a panel of cancer cell lines. Our classification provides an objective ranking for potential CT genes, which is useful in guiding further identification and characterization of these potentially important diagnostic and therapeutic targets.**

gene index | prediction

**C**ancer/Testis (C/T) genes are a heterogeneous group that are normally expressed predominantly in germ cells and in trophoblasts, and yet are aberrantly activated in up to 40% of various types of cancer types (1). A subset of the CT genes has been shown to encode antigens that are immunogenic and elicit humoral and cellular immune responses in cancer patients (2). Because of their restricted expression profile in normal tissues and because the testis is an immunoprivileged site, the CT antigens are emerging as strong candidates for therapeutic cancer vaccines, as revealed by early-phase clinical trials (3–10). Biologically, the CT genes provide a model to better understand complex gene regulation and aberrant gene activation during cancer.

Any gene that exhibits an mRNA expression profile restricted to the testis and neoplastic cells can be termed a CT gene. Existing definitions of CT genes vary in the literature, from genes expressed exclusively in adult testis germ cells and malignant tumors (1, 11) to dominant testicular expression (12), possible additional presence in placenta and ovary and epigenetic regulation (13), or membership of a gene family and localization on the X chromosome (14). Reflecting this lack of a consensus definition, an increasing number of heterogeneous CT candidates have appeared in the literature, with available

expression profile information frequently limited to the original defining articles. In some cases, e.g., ACRBP, the original CT-restricted expression in normal tissues could not be confirmed by subsequent experiments (1). Partially due to this lack of a clear and broadly applicable definition, or “type specimen,” for a CT gene, it has become increasingly challenging to identify the CT genes that are most suitable for cancer vaccine development. Moreover, this incoherent classification increases the risk of pursuing unsuitable clinical targets. However, with more expression data becoming available, CT gene transcripts of genes originally thought to have the CT expression profile are being detected in additional tissues (1), resulting in the more stringent “testis-restricted” description being altered to one of “testis-preference.” Based on a compilation from the published literature, the CT database now lists >130 RefSeq nucleotide identifiers as CT genes that belong to 83 gene families ([www.cta.lncc.br](http://www.cta.lncc.br)). An analysis of the human X chromosome has also suggested that as many as 10% of the genes on this chromosome may be CT genes (15). Given this increasing number of CT and CT-like genes, their comprehensive classification based on expression profiles is essential for our understanding of their biological role and regulation of expression.

In an attempt to resolve this and to identify new CT antigens, we have taken an in silico approach to produce a comprehensive survey of CT gene expression profiles by combining expression information from an existing corpus of >8,000 cDNA libraries (16) together with the depth and resolution provided by massively parallel signature sequencing (MPSS) expression libraries (17), cap-analysis of Gene Expression (CAGE) libraries (18), and a survey using semiquantitative reverse-transcription PCR (RT-PCR) on a panel of 22 normal tissues. As a result, we have created a coherent classification of CT genes, and new CT genes have been identified using well-informed, structured prediction and confirmation criteria.

## Results and Discussion

**CT classification.** CT genes were classified into 3 groups, testis-restricted, testis/brain-restricted and testis-selective, based on

Author contributions: O.H., O.L.C., C.A.M., U.S., A.K., A.J.S., L.J.O., and W.H. designed research; O.H., O.L.C., T.C., and R.C. performed research; B.J.S., Y.-T.C., T.C., C.A.M., S.P., U.S., M.L., A.K., P.C., Y.H., and C.V.J. contributed new reagents/analytic tools; O.H., O.L.C., and B.J.S. analyzed data; and O.H. wrote the paper.

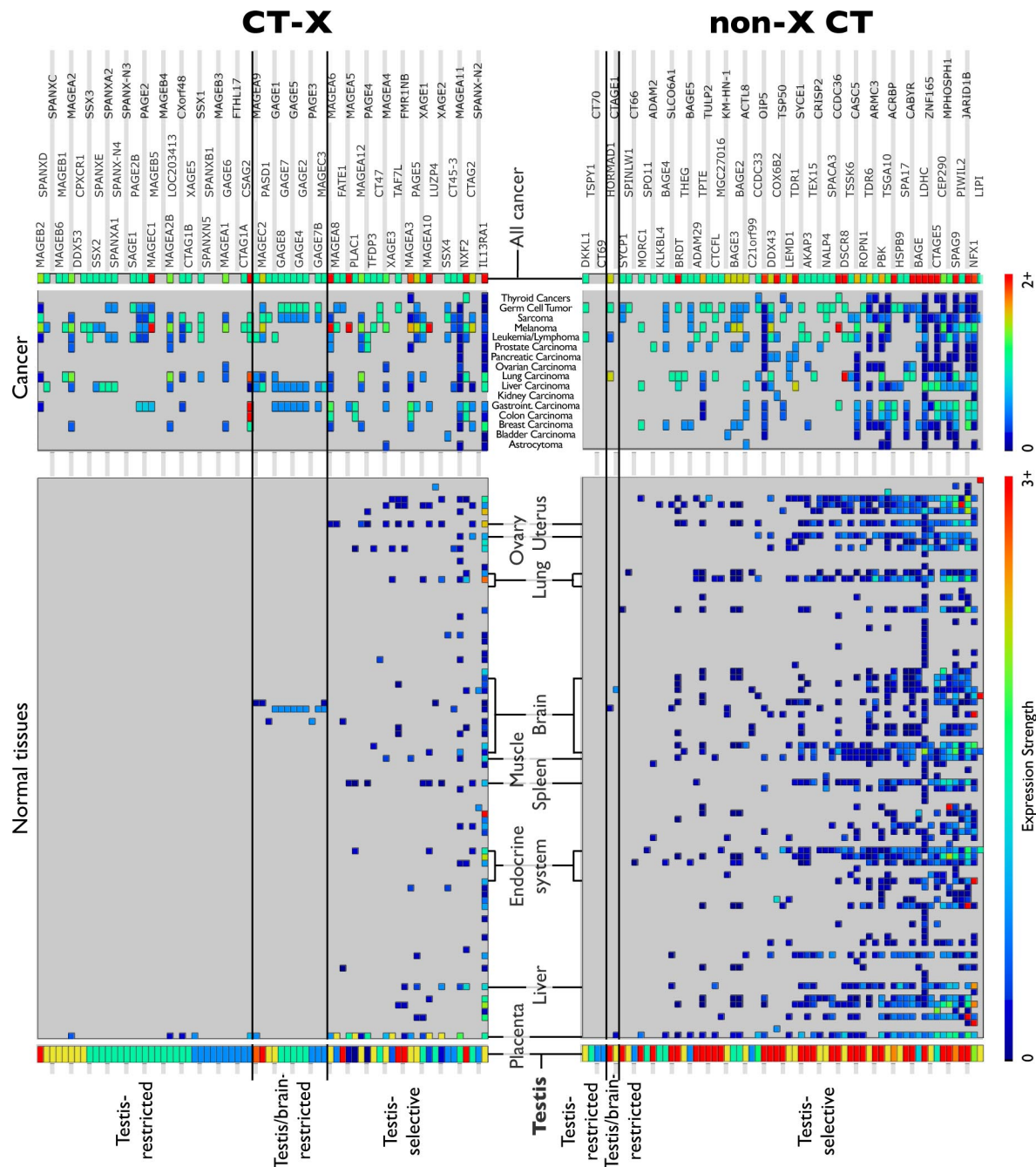
The authors declare no conflict of interest.

Freely available online through the PNAS open access option.

<sup>1</sup>To whom correspondence may be addressed. E-mail: [ohofmann@hsph.harvard.edu](mailto:ohofmann@hsph.harvard.edu) or [lold@licr.org](mailto:lold@licr.org).

This article contains supporting information online at [www.pnas.org/cgi/content/full/0810777105/DCSupplemental](http://www.pnas.org/cgi/content/full/0810777105/DCSupplemental).

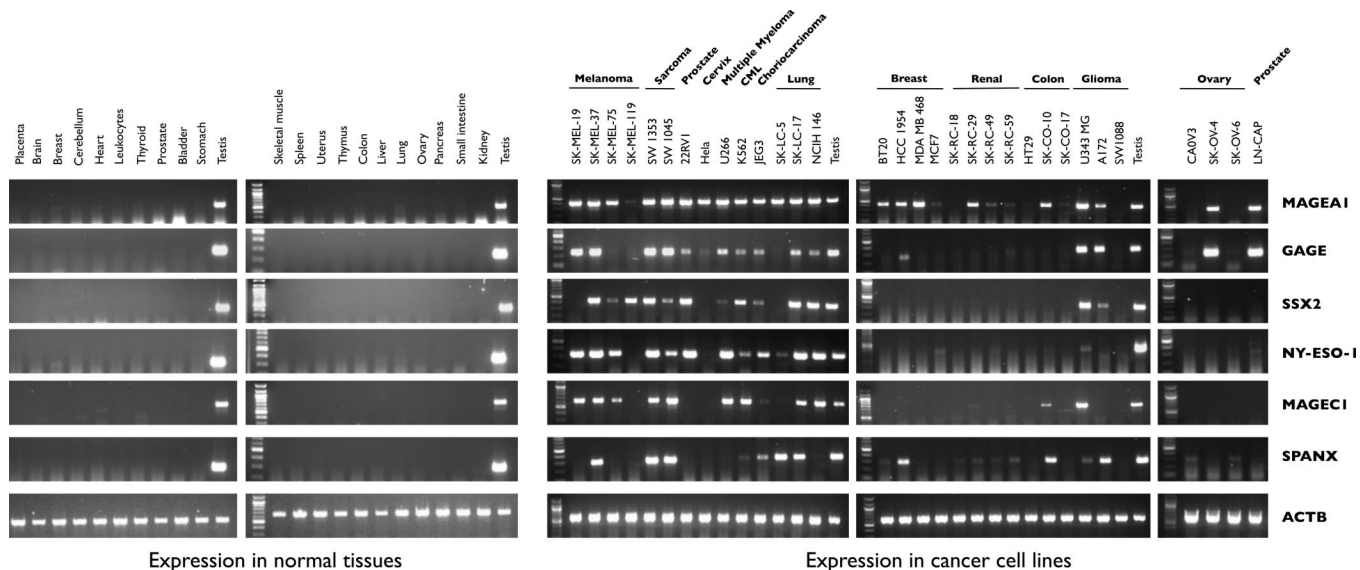
© 2008 by The National Academy of Sciences of the USA



**Fig. 1.** Merged expression profiles of CT-X (left array) and non-X CT genes (Right) based on expression data from RT-PCR and cDNA, MPSS and CAGE libraries from tissues sources annotated as normal and “adult” (Lower) or “cancer.” Expression in normal testis, placenta, and selected tissues is marked. Color reflects the support for the expression of a CT genes in a given anatomical site (blue for low combined expression evidence  $\geq 1$ , red for strong support from at least 3 sources (for the normal tissue panel) with a total score  $\geq 3$ ) or 2 sources (the cancer panel lacking RT-PCR data), respectively. The most abundant expression (red) is seen in testis for most genes, particularly in the non-X CT group. Expression values were normalized on a per-gene basis relative to the combined normal testis/placenta expression confidence (Lower) or the source of the highest cancer expression confidence (Upper). The 3 CT annotation groups (testis-restricted, testis/brain-restricted and testis-selective) are highlighted. See [Dataset S3](#) for the full list of CT classifications.

their expression profiles obtained from a manually curated corpus of cDNA, MPSS, CAGE expression libraries and RT-PCR (see [Dataset S1](#) for MPSS and CAGE library annotation and <http://evocontology.org> for the cDNA annotation). By merging expression information using different technology platforms, we were able to leverage their individual strengths—the breadth of tissue coverage associated with the cDNA/EST expression libraries, the high sensitivity of CAGE/MPSS and the ability to

custom-tailor PCR primers. Of 153 genes, 39 with transcripts present only in adult testis and no other normal adult tissue except for placenta were classified as testis-restricted; 14 CT genes with additional expression in other adult immunorestricted sites (all regions of the brain) were classified as testis/brain-restricted, and 85 genes, designated as testis-selective, were ranked by the ratio of testis/placenta expression relative to other expression in normal adult tissues (see Fig. 1 for



**Fig. 2.** RT-PCR analysis of selected CT genes in the testis-restricted category (MAGEA1, GAGE, SSX2, NY-ESO-1, MAGEC1, and SPANX). Expression profile are shown for a range of 22 normal tissues (Left) and 31 cancer cell lines (Right).

the expression array, Fig. 2 for the PCR panel of selected testis-restricted CT genes, and Fig. S1 and Dataset S2 for arrays from individual expression sources).

An uneven chromosomal distribution of the CT genes was observed, with 83 of 153 genes (54%) being on the X chromosome, and 70 on non-X chromosomes (Fig. S2). Furthermore, 35 CT-X genes were classified as testis-restricted, whereas only 4 non-X CT genes belong to this group. An additional 12 CT-X genes were found to be testis/brain-restricted, compared with 2 non-X testis/brain-restricted CT genes. CT-X gene family members thus appear to be under more stringent transcriptional restriction in somatic tissues, whereas non-X CT genes are more broadly expressed. This validates the CT gene classification into CT-X and CT non-X groups, with the CT-X group being of particular interest for therapeutic approaches.

Twenty-six CT-X and 59 non-X CT genes belong to the testis-selective category, and 36 of these genes (5 CT-X and 31 non-X CT) had >50% of the expression evidence derived from non-testis or placental libraries, indicating that these might not qualify as CT genes.

Seven CT genes were not identified in any library at all (2 CT-X and 5 non-X CT). An additional 8 CT-X genes (SPANX-N1, PAGE1, CSAG1, SSX5/6/7/9, and CT45-2) were not present in any testis-annotated library. Of these, SSX5 and SSX7 have been shown to be expressed in testis by RT-PCR (19), suggesting a likely discrepancy in mapping short sequence tags to their genomic counterparts, an expected phenomenon for large and highly homologous gene families like SSX. In contrast, the absence of testicular expression of SSX6 and SSX9 was confirmed in that study, indicating that some of the currently recognized CT genes could either be silent or expressed at extremely low levels in testis. The full list with classification and raw expression scores across the merged expression array can be found in Dataset S3.

Associations between different CT gene properties and their assigned classification were analyzed using the APRIORI algorithm. Besides being more likely testis-restricted, CT-X genes were found to be more often members of multigene families than non-X CTs. In addition, Gene Ontology terms showed CT-X genes to be more often in the “molecular function unknown” and “biological process unknown” categories, whereas the non-X CTs are associated with known functions such as meiosis, sexual

reproduction, and gametogenesis (see Dataset S4 for all attributes and annotations).

While the description of CT-X genes such as NY-ESO-1 (20), SSX2 (21), and MAGE-A1 (22) match our classification—all are in the testis-restricted category—not all CT genes were found to be as testis-restricted as described in the literature. BAGE, SPO11, LIPI, LDHC, and BRDT, considered to be testis-restricted based on a tissue panel of 13 non-gametogenic normal tissues (1), fall into the testis-selective category in our screen, most likely due to a larger amount of expression sources sampled. Despite the broader coverage we could not confirm an expression of MAGE-A1, MAGE-C1, and NY-ESO-1 at low levels in the pancreas reported in the same study. In agreement with the study in ref. 1, we found IL13RA1, ACRBP, and SPA17 to be expressed in a wide variety of tissues, falling into the lower end of the testis-selective category.

In the present study, we have ranked the testis-selective genes based upon the ratios of their expression evidence in testis and placenta relative to other somatic tissues, rather than using fixed thresholds and the number of somatic tissues in which a CT candidate is allowed as the distinguishing criteria for CT versus non-CT genes (2). Genes without any somatic expression have unique potential for cancer vaccines and other therapeutic approaches to cancer. From past work involving screening of larger sets of genes (23), a cutoff was introduced that defined CT candidate genes as genes with 2-fold higher expression evidence in testis and placenta relative to all other somatic normal tissues. This approach was complementary to our current one and will not require updated thresholds as the number of sampled tissue sources increases.

Intriguingly, a number of CT genes were found to be expressed in no somatic tissues except for brain, suggesting the presence of a distinctive transcriptional control mechanism that functions with tissue specificity in germ cells and in brain. There have been relatively few studies of CT gene expression in different anatomical regions of normal brain and similarly not many in brain tumors (24, 25), except for NXF2, which was shown to be expressed in normal brain (26). Our in silico study has discovered a broader subset of CT genes with brain expression, among them members of the otherwise fully testis-restricted GAGE and MAGE families, found to be expressed in the hippocampus and cerebral cortex. A previous study has similarly identified a group



genes, 83 that encode 107 RefSeq transcripts were mapped to the X chromosome (CT-X genes) whereas 70 genes were on autosomes (non-X CT genes). Subcellular localization was based on predictions in the human version of the LOCATE system (40). SEREX information was obtained from the Cancer Immunome Database website (<http://ludwig-sun5.unil.ch/CancerImmunomeDB>). Ambiguities were resolved by manual curation.

**Source of Expression Information.** Gene expression profiles were determined based on 4 different sources: 99 CAGE libraries from the RIKEN FANTOM3 project (18), 47 MPSS libraries (17, 23, 41), a collection of 8401 cDNA expression libraries from the eVOC system (16), and semiquantitative RT-PCR across 22 normal tissue samples. Source materials were annotated with regards to the anatomical site and pathological status of their source tissues. In cases where the anatomical source was unclassifiable, cell type information was used. Bone marrow/blood libraries were designated bone marrow, and all combinations with mucosa (colon, stomach) were merged into "mucosa." Libraries not explicitly annotated as "normal" were considered as unclassified. Libraries from pooled tissue sources were ignored, and pooled samples were kept as long as the pathological and anatomical status was identical for all donors (see [Dataset S1](#) for annotated libraries).

**Pseudoarrays.** Expression information was organized into "pseudoarrays" based on expression information obtained from CAGE-, MPSS-, and cDNA-libraries in the case of cancer expression and merged with RT-PCR results in the case of normal tissue expression. Columns reflect the class of library in which a CT transcript was identified and rows represent individual RefSeq transcripts. Annotation was based on the general library class description (normal, cancer or unclassified) combined with pathological state and anatomical site. To evaluate the relative levels of CT expression we converted expression signals from the 4 sources into "expression evidence": For CAGE- and MPSS-based expression data, expression evidence was based on detected tags per million (TPM), with matches <3 TPM ( $\approx 1$  transcript per cell) filtered out. Normalized and subtracted EST libraries prevent quantitation of expression strength based on EST counts, therefore expression evidence is represented by the number of cDNA libraries in which a given transcript was identified. RT-PCR results were manually binned into 5 groups of expression, ranging from 0 (not expressed) to 4 (strongly expressed). For each expression source, evidence values were normalized on a per-transcript basis by setting the highest expression evidence in normal tissues to a value of 1, reflecting relative changes in expression levels across tissues and pathological states. Pseudoarrays from the 4 expression sources were merged by summing the individual expression evidence scores for a given transcript from each platform. Expression profiles for multiple transcripts associated with the same gene were merged into a single representation, keeping the highest expression score for overlapping annotations. In arrays where annotation was "merged" into single columns based on their class (e.g., all cancer expression information), the highest expression score across all annotated libraries was kept for each gene.

**Visualization and Ranking.** Genes were divided into CT-X and non-X CT panels, then individually ranked by their expression properties in normal tissues and classified into the following 3 categories: (i) expression in testis and placenta only (testis-restricted); (ii) expression in testis, placenta and brain-regions only (testis/brain-restricted), and (iii) all other genes (testis-selective). Final ranking within each category was obtained by sorting based on decreasing level of normal tissue specificity as measured by the combined testis and placenta expression evidence divided by all normal expression evidence. All arrays were visualized using MeV 4.0 ([www.tm4.org](http://www.tm4.org)).

**Clustering Methods.** Associations between CT annotation and their classification were investigated by recording their assigned class; presence or absence

in placenta, brain, testis, and developing ovary; their testis/placenta tissue specificity; their X vs. non-X chromosomal status; membership in a gene family; subcellular localization; and evolutionary status (36) followed by an analysis with the APRIORI algorithm (42), which identifies association rules matching a predefined threshold of support (30%) and confidence ( $\geq 0.8$ )

**Search Criteria for CT Candidates.** CT candidates were identified using the same in silico expression sources, but with no filters for minimum TPM value and satisfying the following criteria: (i) exhibit expression in testis and at least one cancer-associated tissue at 10 TPM (CAGE, MPSS) or presence in at least one EST/cDNA library with testis and cancer annotation; (ii) not be present above those levels in any other tissue except for placenta, ovary, and brain; and (iii) be supported independently by 2 platforms. Identified candidates were ranked using the same approach used to classify known CT genes. To increase coverage of CT-X genes, a second genome-wide search was conducted requiring support from only a single platform. Candidates were selected for RT-PCR validation by manual curation, removing hypothetical proteins, predicted genes and candidates with multiple publications indicating expression in somatic tissues.

**RT-PCR.** RNA preparations were purchased from the normal tissue panels of Clontech and Ambion or prepared from cancer cell lines using the RNAeasy kit (Qiagen) and were used to prepare cDNA for RT-PCR. A total of 1.0  $\mu$ g of RNA was reverse transcribed into cDNA in a total volume of 20  $\mu$ L using the Omniscript RT kit (Qiagen) according to the manufacturer's protocol using oligo(dT)<sub>18</sub> primers (Invitrogen). The cDNA was diluted 5 times and 3  $\mu$ L was used in the PCR with primers specific to each analyzed gene in a final volume of 25  $\mu$ L. Primers used for PCR amplification were designed to have an annealing temperature  $\approx 60^\circ\text{C}$  using Primer3 software ([www.genome.wi.mit.edu/cgi-bin/primer/primer3www.cgi](http://www.genome.wi.mit.edu/cgi-bin/primer/primer3www.cgi)) and were chosen to encompass introns between exon sequences to avoid amplification of genomic DNA. DNase treatment was undertaken before cDNA synthesis to analyze intronless genes. Primers were designed to target all known variants of a gene in RefSeq and their specificity was confirmed by aligning with the National Center for Biotechnology Information sequence databases using BLAST ([www.ncbi.nlm.nih.gov/blast/blast.cgi](http://www.ncbi.nlm.nih.gov/blast/blast.cgi)). Primer sequences and amplicon sizes are provided in [Dataset S7](#).

JumpStart REDTaq ReadyMix (Sigma Aldrich) was used for amplification according to the manufacturer's instructions. Samples were amplified with a pre-cycling hold at  $95^\circ\text{C}$  for 3 min, followed by 35 specific cycles of denaturation at  $95^\circ\text{C}$  for 15 seconds, annealing for 30 seconds (10 cycles at  $60^\circ\text{C}$ , 10 cycles at  $58^\circ\text{C}$  and 15 cycles at  $56^\circ\text{C}$ ) and extension at  $72^\circ\text{C}$  for 30 seconds followed by a final extension step at  $72^\circ\text{C}$  for 7 min.  $\beta$ -actin was amplified as control. PCR products were separated on 1.5% agarose gels stained with ethidium bromide. For semiquantitative PCR analysis, RT-PCR products were classified into 0 (negative) to 4 (strongest signal) based on the intensity of the product on ethidium bromide-stained gels.

**ACKNOWLEDGMENTS.** We thank Dmitry Kuznetsov for providing access to the SEREX information on CT genes and Erika Ritter (Ludwig Institute for Cancer Research, New York Branch at Memorial Sloan-Kettering Cancer Center, New York) for providing cell lines. This project was supported by the South African National Bioinformatics Network; National Institutes of Health Stanford-South African Informatics Training for Global Health Grant TW-03-008; Atlantic Philanthropies; The Oppenheimer Memorial Trust; a Research Grant for the RIKEN Genome Exploration Research Project from the Ministry of Education, Culture, Sports, Science and Technology of the Japanese Government (to Y. H.); and a grant from the Genome Network Project from the Ministry of Education, Culture, Sports, Science and Technology, Japan. This work was conducted as part of the Hilton-Ludwig Cancer Metastasis Initiative, funded by the Conrad N. Hilton Foundation and the Ludwig Institute for Cancer Research.

- Scanlan MJ, Simpson AJG, Old LJ (2004) The cancer/testis genes: Review, standardization, and commentary. *Cancer Immunol* 4:1.
- Simpson AJG, Caballero OL, Jungbluth A, Chen YT, Old LJ (2005) Cancer/testis antigens, gametogenesis and cancer. *Nat Rev Cancer* 5:615–625.
- Marchand M, et al. (1999) Tumor regressions observed in patients with metastatic melanoma treated with an antigenic peptide encoded by gene MAGE-3 and presented by HLA-A1. *Int J Cancer* 80:219–230.
- Davis ID, et al. (2004) Recombinant NY-ESO-1 protein with iscomatrix adjuvant induces broad integrated antibody and CD4(+) and CD8(+) t cell responses in humans. *Proc Natl Acad Sci USA* 101:10697–10702.
- Jäger E, et al. (2006) Recombinant vaccinia/fowlpox NY-ESO-1 vaccines induce both humoral and cellular NY-ESO-1-specific immune responses in cancer patients. *Proc Natl Acad Sci USA* 103:14453–14458.
- Valmori D, et al. (2007) Vaccination with NY-ESO-1 protein and CPG in montanide induces integrated antibody/th1 responses and CD8 t cells through cross-priming. *Proc Natl Acad Sci USA* 104:8947–8952.
- Uenaka A, et al. (2007) T cell immunomonitoring and tumor responses in patients immunized with a complex of cholesterol-bearing hydrophobized pullulan (chp) and NY-ESO-1 protein. *Cancer Immunol* 7:9.
- Odunsi K, et al. (2007) Vaccination with an NY-ESO-1 peptide of HLA class II/i specificities induces integrated humoral and t cell responses in ovarian cancer. *Proc Natl Acad Sci USA* 104:12837–12842.
- Atanackovic D, et al. (2006) Expression of cancer-testis antigens as possible targets for antigen-specific immunotherapy in head and neck squamous cell carcinoma. *Cancer Biol Ther* 5:1218–1225.
- Gnjatic S, et al. (2006) NY-ESO-1: Review of an immunogenic tumor antigen. *Adv Cancer Res* 95:1–30.
- Scanlan MJ, et al. (2002) Identification of cancer/testis genes by database mining and mRNA expression analysis. *Int J Cancer* 98:485–492.
- Zendman AJW, Ruiter DJ, Muijen GNPV (2003) Cancer/testis-associated genes: Identification, expression profile, and putative function. *J Cell Physiol* 194:272–288.

