

Literature LabTM analysis of pathway associations relevant to prognostic gene signatures for breast, bone and lung cancer metastasis

Damon Anderson, PhD, Applications Scientist and Paul Martinez, President and CEO, Acumenta Biotech, Westminster, MA. Co-responding author: Damon Anderson, PhD, danderson@acumenta.com, 412-901-7785

The landmark breast cancer metastasis study by van't Veer, et al. 2002. Nature (415) 530-536 identified a 70-gene signature derived from microarray analysis of 117 patients that was strongly predictive of a short interval to distant metastases ('poor prognosis' signature) in patients without tumor cells in local lymph nodes at diagnosis (lymph node negative). We applied Literature Lab[™] analysis on this gene signature to identify significantly associated pathways, thus offering a modern perspective based on the current PubMed content and Literature Lab[™] analytics. Published gene signatures derived from microarray data from lung and bone metastatic tumors were then analyzed using Literature Lab[™] and significant pathways were identified and compared. The Id signaling pathway showed a strong association with lung while the CXCR4 pathway was strongly associated with both lung and bone, results well supported in the literature. Importantly, Literature LabTM identified the IL1 pathway, which was exclusively associated with lung, and the CCR5 pathway. which was significantly associated with lung and bone. The unexpected associations of these immune factor pathways and their roles in the lung and bone metastatic microenvironments indicate potential novel areas for future investigations. Moreover, these results and others showcase the powerful capabilities of Literature Lab[™] as a genetic analysis platform that can offer validation and uncover meritorious functional associations missed by other genetic analysis tools.

| Click on column header | s or buttons to s | ort columns | or rows. Click ag | Diseas Based on gain to reverse | es-MeSH/I PubMed abstract | Pathways s from 01/01/1 ter the first few | - | | | L | i di 29 (?) N T 7 | | |
|-------------------------------|-------------------|---------------|-------------------|---------------------------------------|------------------------------|---|----------------|-------------------------|---------|---|----------------------|-------------------|----------------------------|
| Cens contain the LPP for | r the terms. Clic | K ON A CEN LO | Visualizer | | IIS. RIGHT CIICK | Close this | Subset | | Row Ter | m: | Column T | erm: | Find |
| Term | Abstract Count | View | Breast Cancer | ERBB2 | ErbB Signaling | Estrogen Signaling | Antigen | E-cadherin Signaling | BRCA1 | mRNA Surveillance | EGF | RNA Polymerase | Basal Cell Carcinoma |
| Abstract Count | | 301103 | 39669 | 8930 | 9937 | 45592 | 310335 | 7566 | 6634 | 4555 | 37456 | 92157 | 3307 |
| View Genes | | | | | | | | | | | | | |
| Carcinoma, Ductal, Breast 🕨 | 12261 | | -2.4598 | -3.0156 | -3.0919 | -3.3238 | -4.0610 | -4.2173 | -4.3255 | -4.4215 | -4.4965 | -4.5620 | -4.6252 |
| Carcinoma, Non-Small-Cell Lur | 32779 | | -4.7004 | -3.8775 | -3.6533 | -5.3683 | -4.0492 | -4.2147 | -4.9749 | <u>-3,7391</u> | -2.8713 | -4.0073 | -4.2074 |
| Disease Models, Animal | 275869 | | -3.4349 | -3.9188 | -3.9293 | -3.7172 | -2.3996 | -4.3978 | -4.9954 | -4.5236 | -3.8653 | -3.1677 | -4.8540 |
| Bronchial Neoplasms | 39034 | | -4.7393 | -3.9359 | -3.7100 | -5.4122 | -4.0738 | -4.2559 | -5.0328 | <u>-3.7890</u> | -2.9324 | -4.0099 | -4.1928 |
| Inflammatory Breast Neoplasn | 215 | | -3.5684 | -4.2005 | -4.4212 | -4.9085 | -5.2222 | -4.6550 | - | -5.9909 | -5.2158 | -6.3427 | |
| Breast Neoplasms, Male | 2304 | | -4.0240 | -4.8524 | -4.8988 | -4.9584 | -5.8716 | -7.2413 | -4.4238 | - | -6.0275 | -7.3728 | -6.8819 |
| Small Cell Lung Carcinoma | 6083 | | -5.0022 | -5.1775 | -4.7986 | -5.7582 | -4.4801 | -4.9781 | -6.4018 | <u>-4.7192</u> | -4.0470 | -4.5009 | -4.2406 |
| Hereditary Breast and Ovariar | <u>144</u> | | -4.5289 | -6.1092 | -6.1556 | -5.6131 | -6.6960 | - | -3,4226 | - | -6.7319 | - | - |
| Pulmonary Sclerosing Hemang | 146 | | | <u>-6.1152</u> | <u>-6.1616</u> | <u>-6.8232</u> | <u>-6.4521</u> | - | - | - | <u>-6.7379</u> | -7.1289 | <u></u> |
| Palatal Neoplasms | <u>1180</u> | | - | -6.4207 | -6.4671 | l. | -7.1658 | - | 1 | - | -7.6454 | -8.0364 | -6.5913 |
| Nose Neoplasms | 8588 | | | -7.2827 | -7.3291 | - | -5.9450 | -6.6086 | -7.7557 | <u>-6.6381</u> | <u>-6.8172</u> | -6.4902 | <u>-6.0554</u> |
| Skull Base Neoplasms | 2526 | | - | -7.3533 | -7.3997 | - | -7.4963 | -7.2813 | - | 14 - 14 - 14 - 14 - 14 - 14 - 14 - 14 - | -7.9760 | -8.3670 | <u> </u> |
| Solitary Pulmonary Nodule | 2690 | | | -7.3806 | -7.4270 | | -7.7175 | - | | - | -8.0033 | -8.3943 | <u></u> |
| Mandibular Neoplasms | <u>4379</u> | | -8.2398 | -7.5922 | -7.6386 | - | -7.9291 | - | - | - | - | -8.0038 | -5.7629 |
| Spinal Neoplasms | 7362 | | -6.3071 | - | 8- | -8.5259 | -7,4503 | - | | - | -7.8385 | -8.8315 | -7.3864 |

The Literature LabTM basic viewer is shown above with co-occurrences tabulated between diseases (breast, lung, bone neoplasms; y-axis) and pathways (x-axis). The pathways are ranked by the strength of their associations in the literature with Carcinoma, Ductal, Breast from highest strength association according to LPF (log of product of frequency). Blue dots hyperlink to the genes associated in the literature with the disease or pathway. All other data is hyperlinked directly to the PubMed literature.



The log of product of frequency (LPF) is a quantitative measure of the strength of association based on the fractional co-occurrence between a term/gene and a term. The Literature LabTM basic LPF is defined as $Log(X/T1 \times X/T2)$, where X/T1 is the percentage of abstracts that mention the first term and X/T2 is the percentage of abstracts that mention the second term. The Literature LabTM PLUS LPF is defined as $Log(X/G \times X/T)$, where X/G is the percentage of gene abstracts that mention the term and X/T is the percentage of term abstracts that mention the given gene. The LPF is not sensitive to abstract volume and the closer to zero the stronger the association.

Literature Lab[™] PLUS was employed to analyze a 70-gene signature derived by DNA microarray analysis on primary breast tumors of 117 young patients. This gene expression signature is strongly predictive of a short interval to distant metastases ('poor prognosis' signature) in patients without tumor cells in local lymph nodes at diagnosis (lymph node negative).

Gene expression profiling predicts clinical outcome of breast cancer

Laura J. van 't Veer, et al. *Nature***415**, 530-536 (31 January 2002)

Breast cancer patients with the same stage of disease can have markedly different treatment responses and overall outcome. The strongest predictors for metastases (for example, lymph node status and histological grade) fail to classify accurately breast tumours according to their clinical behaviour. Chemotherapy or hormonal therapy reduces the risk of distant metastases by approximately one-third; however, 70–80% of patients receiving this treatment would have survived without it. None of the signatures of breast cancer gene expression reported to date allow for patient-tailored therapy strategies. Here we used DNA microarray analysis on primary breast tumours of 117 young patients, and applied supervised classification to identify a gene expression signature strongly predictive of a short interval to distant metastases ('poor prognosis' signature) in patients without tumour cells in local lymph nodes at diagnosis (lymph node negative). In addition, we established a signature that identifies tumours of *BRCA1* carriers. The poor prognosis signature consists of genes regulating cell cycle, invasion, metastasis and angiogenesis. This gene expression profile will outperform all currently used clinical parameters in predicting disease outcome. Our findings provide a strategy to select patients who would benefit from adjuvant therapy.

Statistically significant associations are scored according to qualifiers as shown below:

| Tier | | | | |
|----------|----------|--------------|-----------------|-----------------------|
| >1024 | Positive | Moderate | Strong | h Parameters |
| 513-1024 | 1.00 | 1.50 | 2.00 | F Score >= |
| 257-512 | 0.1587 | 0.0668 | 0.0228 | Value <= |
| 129-256 | 75 | 50 | 25 | rm Rank <= |
| 65-128 | 10.0 | 15.0 | 20.0 | NonZero Genes >= |
| 33-64 | 2 | 2 | 3 | te Min NonZero Genes |
| 17-32 | 10 | 25 | 50 | X Term Abstracts >= |
| 9-16 | 100.0 | 99.0 | 98.0 | Single Gene % Contrib |
| 5-8 | 250 | 500 | 750 | Random Sets >= |
| 3-4 | | - | | |
| 2 | | |) Default Value | Set t |
| 1 | | lues for use | aca Qualifier v | You may save th |
| 0 | | ides for use | ese Qualmer vi | rou may save u |
| Total | | view by | ments that yo | with all exper |



The P-value and other measures of statistical significance are used to derive the strength of association. The heuristics involved reward associations for representation of more genes in a set and downgrade associations driven by sparse data. The qualifier criteria can be adjusted according to experimental goals but typically default values are used for straight forward analysis. Importantly, Literature LabTM is the only platform that treats each gene list (and gene) as unique and finds significant associations in literature based on LPF and statistical qualifiers. Other methods search for associations between an experimentally derived gene list and pathway or domain information within curated databases, without regard for the uniqueness of each gene set.

The 'Highlights' tab displays the MeSH term domains that were used in the analysis (left column of left panel) and the number and strength of associations that were identified based on the qualifiers (right side of left panel). The right panel displays the identity of the associations, the score, and the p-value. Strong (green), Moderate (blue), and Positive (pink) pathway associations are shown:

| | Strong | Moderat | e Positive | Visualize |
|-----------------------|-----------|--------------------|-------------|-----------|
| Pathways | 9 4 | <u>10</u> | 23 | |
| Diseases-MeSH | 01 | 9 2 | <u>9 45</u> | |
| PathologicalCond-MeSH | 9 5 | <u>5</u> | 9 2 | - |
| PharmacoAction-MeSH | 0 | <u>1</u> | 9 9 | |
| PharmacoSubst-MeSH | 2 | <u>) 3</u> | 23 | |
| ChemicalActions-MeSH | 2 | 4 | 0 2 | - |
| ChemicalsDrugs-MeSH | 15 | 4 6 | 179 | - |
| Anatomy-MeSH | 2 | 0 6 | 24 | |
| Physiology-MeSH | <u>3</u> | 9 <u>12</u> | 16 | - |
| CellPhysiology-MeSH | <u>8</u> | 🥥 Z | 🥥 Z | |
| CellTypes-MeSH | 0 | <u> 5</u> | <u>8</u> | - |
| CellStructures-MeSH | 🔵 Z | 🥥 Z | 🥥 <u>4</u> | |
| <u>CellLines</u> | 2 | 5 | 10 | - |
| TissueTypes-MeSH | 01 | 5 | 0 11 | |
| Metabolism-MeSH | 2 | 0 | 3 | - |
| Biogenetics-MeSH | 4 | <u>16</u> | 🥥 Z | |
| OtherBiology-MeSH | <u>5</u> | 1 7 | <u>) 17</u> | - |
| Organisms-MeSH | 0 | 01 | 🥥 11 | |
| Substances-MeSH | 3 | <u>8</u> | <u>) 34</u> | - |
| Psychology-MeSH | 0 | 01 | 0 | - |

Associations include a number of pathways well documented in cancer growth and development including: Cell Cycle/Cyclins, DNA replication/G1/S Phase, etc. Other pathway associations have well documented ties to metastasis including: Angiogenesis/VEGF, Matrix Metalloproteinase, Beta-Arrestin, etc. Pathway associations representing interesting areas of investigation in the literature are also uncovered. For instance, the Insulin pathway has been the target of research and therapeutic discovery over the past few years. However, the complex cross talk between factors and receptors has made it elusive in the treatment of breast carcinoma.

VEGF 3.57 0.0002 3.16 0.0008 Angiogenesis 2.27 0.0116 p53 **DNA Replication** 2.03 0.0212 Moderate Assocations for Pathways Term P-Value Score G1/S Checkpoint 2.95 0.0016 2.76 0.0029 Cyclin 2.51 0.0060 Cell Cycle IGF-1 2.49 0.0063 UPAR 2.35 0.0095 IGF-1R 2.10 0.0180 TGF Beta 2.07 0.0194 CDC42 1.84 0.0332 0.0360 1.80 FOXM1 Ubiguitin Mediated Proteolysis 1.80 0.0362

Strong Assocations for Pathways

| Term | Score ¥ | P-¥alue | |
|--------------------------|---------|---------|--|
| Matrix Metalloproteinase | 4.31 | 0.0000 | |
| Gemcitabine | 3.61 | 0.0002 | |
| p27 Phosphorylation | 2.68 | 0.0036 | |
| <u>CDK</u> | 2.59 | 0.0048 | |
| RAS RHO | 2.24 | 0.0126 | |
| Retinoblastoma | 2.22 | 0.0133 | |
| RB Tumor | 2.07 | 0.0193 | |
| <u>p73</u> | 2.05 | 0.0203 | |
| E2F | 2.03 | 0.0211 | |
| TGF-beta Receptor | 2.00 | 0.0228 | |
| PI3K | 1.85 | 0.0319 | |
| Antisense | 1.85 | 0.0320 | |
| FGF | 1.76 | 0.0389 | |
| G2/M Checkpoint | 1.76 | 0.0391 | |
| RNA Polymerase | 1.76 | 0.0392 | |
| Aurora Kinase | 1.70 | 0.0443 | |
| Insulin | 1.62 | 0.0525 | |
| Breast Cancer | 1.62 | 0.0525 | |
| Visual Signal | 1.61 | 0.0537 | |
| Hypoxia p53 | 1.51 | 0.0661 | |
| Beta-Arrestin | 1.48 | 0.0693 | |
| C-MYB | 1.44 | 0.0755 | |
| TSP-1 | 1.18 | 0.1182 | |

| Click on column head Percentages = prop | lers or buttons | to sor ment ! | Ex t columns or Set LPF contr | cperiment: r rows. Click aga ributed by a gen | Copy of v ain to reverse ne. Click on a | ant Veer M Based on Publ sort order. Ent cell to view the | Domain letastasizi led abstracts fro ter the first few abstracts for t | n: Pathway ng Breast m 01/01/1990 r letters of a ge he gene and te | /S Cancer Sig through 02/28/2 ene and/or tern erm. Right Click | nature Se 2015 In to find a cell. | pt 2015 up 15. | odate 2 | U | o 4 3 N T . |
|--|-----------------|------------------|-------------------------------------|--|---|--|--|--|--|---|-------------------|----------|--------------------|----------------|
| <u>.</u> | Clear/Set/Use C | heck Ma | arks | | Create a Subs | et of the Domain | | Visu | alizer 😽 | | Gene: | Te | erm: | Find |
| Pathways | Count | ~ | VEGF | Angiogenesis | p53 | DNA Replication | G1/S Checkpoint | Cyclin | Cell Cycle | IGF-1 | UPAR | IGF-1R | TGF Beta | CDC42 |
| iew All Genes For Term | | | | | | | | | | | | | 1000 * 1000 | |
| sociation 🕨 | | | Strong | Strong | Strong | Strong | Moderate | Moderate | Moderate | Moderate | Moderate | Moderate | Moderate | Moderate |
| ne X Term Abstracts | | 2.31 | 3769 | 2761 | 1324 | 1029 | 1786 | 3426 | 5432 | 550 | 514 | 152 | 2145 | 12 |
| rm Abstracts | | 2 3 | 27043 | 35926 | 41025 | 31771 | 18198 | 27450 | 170719 | 15523 | <u>5004</u> | 4511 | 42664 | 463 |
| nzero Genes | | - | 27 | 25 | 37 | 32 | 26 | 30 | 50 | 17 | 9 | 11 | 33 | |
| F | | | -1.29 | -1.89 | -2.08 | -2.29 | -1.56 | -1.15 | -1.87 | -2.19 | -2.66 | -2.83 | -1.90 | -2.2 |
| ndom Sets 🛨 | Count | | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 1000 | 100 |
| periment Set | Term rank | | 10 | 10 | 14 | 17 | 8 | 21 | 26 | 44 | 27 | 43 | 44 | |
| | Score | | 3.57 | 3.16 | 2.27 | 2.03 | 2.95 | 2.76 | 2.51 | 2,49 | 2.35 | 2.10 | 2.07 | 1.8 |
| | P-Value | | 0.0002 | 0.0008 | 0.0116 | 0.0212 | 0.0016 | 0.0029 | 0.0060 | 0.0063 | 0.0095 | 0.0180 | 0.0194 | 0.033 |
| | Score rank | | 3 | 4 | 20 | 31 | 6 | 8 | 13 | 15 | 17 | 25 | 28 | 2 |
| mbol/OrigID | Max PF | | 95.51% | 83.17% | 79.13% | 64.09% | 98.17% | 98.85% | 84.79% | 97.56% | 98.86% | 97.58% | 93.20% | 89.94 |
| 199 | 21389 | | 3./1% | 13.1/% 0 | 0.4/% | 0.05% | 0.1/% | 0.10% | 1.35% | 0.21% | 98.86%0 | 0.76% | 5.15% | 0.15 |
| FP2 | 5059 | (mm) | 95.51% | 83.1/%0 | 0.07% | 0.01% | 0.00% | 0.00% | 0.65% | 0.06% | 0.80% | 0.14% | 02.205 | 0.01 |
| IFD0 M1 | 2297 | | 0.05% | 0.14% | 0.05% | 0.01% | 0.01% | 0.01% | 0.020/ | 0.59% | 0.02% | 0.16% | 93.20% | 0.07 |
| NE2 | 160 | | 0.02% | 0.119/ | 10 169/ 8 | 64.00% | 09 179/ 1 | 0.01% | 94 708/ 8 | 0.270/ | 0.06% | 0 200/ | 0.00% | 0.03 |
| T203 | 9766 | | 0.03% | 0.02% | 19.10% | 0.00% | 30.1/% | 0.00% | 0.020/ | 0.0/% | 0.10% | 0.38% | 0.00% | 0.03 |
| LEMU INI | 900 | | 0.03% | 1 2094 | 0.03% | 0.00% | 0.00% | 0.00% | 0.03% | 0.00% | 0.01% | 0.03% | 0.00% | |
| SD1 | 186 | | 0.01% | 0.67% | 0.02% | 0.01% | 0.00% | 0.00% | 0.03% | 0.01% | 0.05% | 0.03% | 0.00% | 0.05 |
| E18 | 198 | | 0.03% | 0.11% | 0.00% | 0.0176 | 0.01% | 0.00% | 0.00% | 0.13% | 0.0076 | 0.07% | 0.11% | 0.00 |
| FBP5 | 1205 | | 0.01% | 0.01% | 0.02% | 0.01% | 0.01% | 0.00% | 0.21% | 97 56% | 0.03% | 97.58% | 0.22% | 0.05 |
| C80 | 314 | | 0.0176 | 0.0176 | 0.01% | 0.03% | 0.02% | 0.01% | 1 54% | 57.5076 | 0.0376 | 37.30 /6 | 0.00% | 0.01 |
| K | 753 | | 0.00% | 0.00% | 0.04% | 0.36% | 0.02% | 0.00% | 0.29% | | | | 0.00% | 0.01 |
| N4DI 1 | 1014 | | 0.00% | 0.00% | 0.0478 | 0.30 /8 | 0.0276 | 0.00 /8 | 0.01% | | | | 0.01% | |
| NIDA | 807 | | 0.00 % | 0.00 % | 0.04% | 0.59% | 0.11% | 0.01% | 1.46% | | | | 0.0176 | 0.04 |
| C3 | 614 | | 0.00% | 0.02% | 79 13% | 0.01% | 0.01% | 0.12% | 0.85% | 0.03% | | 0.07% | 0.02% | 0.04 |
| 14 M | 014 | | 0.00 78 | 0.0276 | 13.13/0 | 0.0170 | 0.0176 | 0.1270 | 0.00 /0 | 0.0076 | | 0.0270 | 0.0270 | |

Shown above are the pathway associations in the 'Details' tab of the Literature LabTM PLUS viewer. Down the y-axis are the qualifiers used to establish the strength of associations including: the number of gene abstracts, the number of term abstracts, the number of non-zero genes, the log of product of frequency, the number of random sets, the term rank, the score, the p-value, etc. Also shown is the list of genes and their relative contribution to the LPF with hyperlinks directly to the literature. Across the x-axis are the associations, which can be ordered by any of the rows on the left column. (Note: shown is a truncated list of genes and pathway associations).

The paper discussed functional annotation for the genes in the 70-gene signature, which provided insight into the underlying biological mechanism leading to rapid metastases. Genes involved in cell cycle, invasion and metastasis, angiogenesis, and signal transduction are significantly upregulated in the poor prognosis signature (for example cyclin E2, *MCM6*, metalloproteinases *MMP9* and *MP1*, *RAB6B*, *PK428*, *ESM1*, and the VEGF receptor *FLT1*.

Literature LabTM PLUS found a number of associations between the signature and these gene pathways (*e.g.* VEGF, Metalloproteinases, etc.). Interestingly, some pathways are not significantly associated as indicated by the qualifier scores (*e.g.* MCM6, RAB6B). In the paper, the authors do not detail a comprehensive annotation of all of the genes in the signature. A powerful aspect of Literature LabTM PLUS is that it treats each gene in the signature as unique and derives a statistically significant association of the entire signature based on scoring parameters. Another important aspect to note is the difference in dates of this published work (1/2002) and the current build of the Literature LabTM database (3/2015). Undoubtedly, there has been a progression in published papers following the role of certain genes in specific pathways, thereby leading to an increase in the strength of association during this time. Likewise, other gene/pathway associations have not been as fruitful and therefore the strength has likely diminished.



Clustering analysis on the significant terms and the genes in the70-gene signature was then explored using the Literature LabTM PLUS 'Term/Gene Cluster' tab. Eight clusters of genes with related function were identified and three of these and their associated heat maps are shown below.



Clusters 2 and 3 are related to cell cycle associated pathways including G1 phase, cyclins, mitosis, microtubules, as well as signaling pathways well documented with tumor proliferation including, aurora kinase, retinoblastoma, E2F, and C-MYB, etc. Cluster 6 is related to pathways associated with tumor metastasis including, matrix metalloproteases, CDC42, etc. The heat maps are an indicator of strength of association of each clustered gene with a pathway, and each box is hyperlinked directly to the PubMed literature.

Nature 415, 530-536 (31 January 2002) doi:10.1038/415530a



The landmark paper authored by van't Veer *et. al.* defined a prognostic 70-gene signature that identifies the metastatic potential of breast tumor cells. This signature was originally associated with several cell growth, metastatic, and angiogenesis pathways.

We next wanted to compare pathway associations between this 70-gene metastatic signature and two signatures linked with metastasis targeted to two organ systems: 1) lung and 2) bone. The goal was to identify pathways that are related to the distinction between basic metastatic potential and metastasis targeted to specific tissues, *i.e.* lung and bone.

In the paper entitled **Genes that mediate breast cancer metastasis to lung** (Andy J. Minn, *et al. Nature* **436**, 518-524. July 28 2005), the authors define a 95-gene signature indicative of aggressive lung metastatic behavior by means of transcriptomic microarray analysis of highly and weakly lung-metastatic cell populations. The genes in this signature were largely distinct from those identified in bone metastatic isolates derived from the same parental line (see below). A 54-gene subset of this signature, which was more similar to the lung metastatic populations selected *in vivo* and postulated to serve specialized lung-related function, was selected as a refined signature. Included are genes associated with: EGF and HER/ErbB signaling, MMP1 and MMP2 matrix metalloproteases, and the IL13Ra2 and VCAM1 receptors, the latter indicating specific roles in the lung tumor microenvironment.

In the paper entitled **A multigenic program mediating breast cancer metastasis to bone** (Yibin Kang, *et al. Cancer Cell* **3**, 537-549. June 2003), the authors define a 102-gene signature indicative of osteolytic bone metastasis by means of microarray analysis of highly and weakly bone-metastatic cell populations. Most of these genes encode osteolytic and angiogenic factors and interestingly, none of the genes in the signature overlap with those in the van't Veer signature.

Using the gene list comparison tool, Literature Lab[™] PLUS identified several pathways that are associated with all three of the gene lists, including those involved in angiogenesis, cell adhesion, and apoptosis.





We then looked at specific pathways that were identified in the previous papers as uniquely indicative of either lung or bone metastasis. Previous work identified gene combinations that act synergistically to promote lung metastasis including: ID1, SPARC, IL13Ra2, VCAM1, MMP2, MMP1, CXCL1, EREG, and COX2.Whereas, genes that were identified as promoting bone metastasis include: IL11, OPN, CXCR4 and CTGF.



The Id Signaling pathway (ID1, inhibitor of DNA binding), previously identified as being significantly upregulated in lung metastasis, shows a predictably strong association with the lung metastasis gene set here. Interestingly, while IL1 is an inflammatory cytokine expressed by macrophages typically associated with infection, it shows a moderate association exclusively with the lung metastasis gene set. "Educated" macrophages often arise in response to the tumor microenvironment, and the exclusive moderate association between IL1 and lung metastasis may indicate a promising potential area of investigation.

Recent research has also pointed to the critical role that CXCR4 receptor and its ligand CXCL12 play in the metastasis of various types of cancer. Lung, bones and lymph nodes all secrete high levels of CXCL12, which acts as a chemoattractant that drives CXCR4-positive primary breast tumor cells towards these secondary metastatic sites. Therefore, it is not unexpected that the lung and bone gene lists show strong associations with CXCR4. It is however interesting that CCR5 shows a moderate association with these gene lists.CCR5 is predominantly expressed on T cells, macrophages, dendritic cells, eosinophils and microglia, and it is likely that CCR5 plays a role in inflammatory responses to infection, though its exact role in normal immune function is unclear. CCR5 may be another emerging player in the tumor metastasis microenvironment area that warrants further investigation.

PLK3 (Polo like kinase 3) is a cytokine inducible serine/threonine kinase that has been implicated in a number of cancers including breast cancer. The fact that it is strongly associated with the bone metastasis gene set, to a significantly greater degree than the other two gene sets, implies that there is a focus on the role of PLK3 in bone metastatic cancer. Likewise, PLK3 may be yet another promising target for further investigation.



We next used the Literature LabTM gene set comparison feature to compare the lung and the bone gene signatures and to identify pathways that are significantly different between the two.ID Signaling is moderately different according to the statistical data and qualifiers, consistent with our earlier result. Interestingly, there are a number of pathway differences (*e.g.* Fatty Acid Biosynthesis, Multi-Drug Resistance, Colorectal Cancer, Nitric Oxide, etc.), which may be useful in further characterizing these metastatic subtypes and helping to drive hypotheses and future studies.

| ×Liter | ature × | | | | Experiment 1 E | Pathway) Gene Lists E (periment 2) I | s Domain (Breast Cancer Breast Cancer | Compariso Metastasis Lu Metastasis Bo | n Ing Nature 20 one 3 | 05 | | | L | oi 4 () N T 7 |
|-------------------------|-------------------|------------------------|---------------------|--------------------------|---------------------------------------|---|--|---|-----------------------------|-----------------------------------|---------------|---------------|--------------------|------------------|
| Click on column hea | ders or buttons | to sort | columns or | rows. Click ag | ain to reverse | sort order. En | ter the first few | letters of a ge | ene and/or ter | m to find a cell. | | | | |
| Percentages = prop | portion of Experi | iment Se | et LPF conti | ributed by a ge | ne. Click on a | cell to view the | abstracts for t | he gene and te | erm. Right Clic | k for Annotatio | ns. | | | |
| | | | Show | Differences (cli | ick on Associat | ion row to reve | erse) 🥘 | Show Similari | ties | Show Significa | | | | |
| | Clear/Set/Use C | heck Mar | ks | | Create a Subset of the Domain | | Visualizer 😽 | | | Gene: | | | erm: | Find |
| Pathways | Count | ✓ Fatt Acid Bios | :y I ynthesis | Multi-Drug Resistance | Colorectal Cancer | Id Signaling | Granulocyte Adhesion | Nuclear Receptors | Nitric Oxide | Acute Inflammatory Response | Leukotriene | Cyclophosp | Cytochrome P450 | AHR |
| View All Genes For Term | | 1000 | | | · · · · · · · · · · · · · · · · · · · | | | | | | | | | |
| Association 🕨 | | St | rong 1^ | Strong 1^ | Strong 1^ | Moderate 1 | Moderate 1 | Moderate 1 | Moderate 1 | Moderate 1 | Moderate 1 | Moderate 1 | Moderate 1 | Moderate 1 |
| Gene X Term Abstracts | | | 10297/723 | 220/54 | 1774/210 | 1232/195 | 673/74 | 333/51 | 4289/349 | 679/77 | 515/33 | 89/26 | 475/45 | 519/115 |
| Term Abstracts | | 191 | 915/191915 | 9294/9294 | 27817/27817 | 14864/14864 | 5979/5979 | 10379/10379 | 63726/63726 | 8466/8466 | 6269/6269 | 3393/3393 | 13727/13727 | 11348/11348 |
| Nonzero Genes | | | 36/27 | 25/10 | 37/23 | 28/17 | 16/10 | 27/13 | 26/18 | 20/9 | 13/9 | 14/6 | 22/12 | 33/17 |
| LPF | | | -1.92/-3.90 | -3.33/-4.70 | -2.74/-4.09 | -1.70/-3.90 | -2.37/-4.26 | -2.72/-4.74 | -2.26/-3.88 | -2.98/-4.38 | -3.09/-5.07 | -3.73/-4.83 | -2.67/-5.03 | -2.41/-4.03 |
| Random Sets 🛨 | Count | | 1000/1000 | 1000/1000 | 1000/1000 | 1000/1000 | 1000/1000 | 1000/1000 | 1000/1000 | 1000/1000 | 1000/1000 | 1000/1000 | 1000/1000 | 1000/1000 |
| Experiment Set | Term rank | 1 | 6/796 | 60/758 | 97/855 | 13/699 | 54/476 | 215/917 | 47/526 | 68/579 | 82/700 | 22/438 | 202/788 | 67/475 |
| | Score | 1 | 2.17/-0.90 | 1.47/-0.68 | 1.23/-0.91 | 3.08/-0.59 | 1.94/-0.08 | 0.83/-1.16 | 1.76/-0.22 | 1.65/-0.28 | 1.32/-0.56 | 1.98/0.10 | 1.04/-0.82 | 1.71/-0.14 |
| | P-Value | 0.0 | 0149/0.1849 | 0.0703/0.2469 | 0.1096/0.1826 | 0.0010/0.2763 | 0.0264/0.4671 | 0.2020/0.1237 | 0.0392/0.4135 | 0.0497/0.3912 | 0.0939/0.2862 | 0.0238/0.4595 | 0.1500/0.2074 | 0.0437/0.4442 |
| | Score rank | 1. | 24/576 | 63/537 | 87/579 | 7/509 | 33/334 | 141/608 | 43/372 | 52/397 | 80/499 | 32/277 | 105/569 | 48/351 |
| Symbol/OrigID | Max PF | 94 | .29/37.99% | 88.25/51.70% | 76.19/29.97% | 98.24/26.34% | 86.87/60.12% | 91.49/45.23% | 91.85/49.69% | 45.35/46.01% | 84.35/27.42% | 85.16/47.17% | 93.23/21.13% | 73.53/72.87% |
| PTG52 | 32308 | | 94.29% | 2.93% | 76.19% | 0.89% | 0.17% | 1.89% | 91.85% | 26.86% | 84.35% | 1.95% | 2.40% | 0.94% |
| MMP2 | 19379 | | 0.43% | 0.74% | 8.10% | 0.17% | 0.25% | 0.15% | 1.49% | 2.32% | 0.53% | 0.52% | 0.02% | 0.22% |
| VCAM1 | 9579 | | 0.53% | 0.09% | 0.23% | 0.10% | 86.87% | 0.33% | 4.34% | 8.94% m | 2.23% | 0.41% | 0.03% | 0.11% |
| CXCR4 | 9463/9463 | | 0.04/4.19% | 2.19/51.70% | 1.34/29.97% | 0.02/3.22% | 0.76/60.12% | 0.04/4.54% | 0.11/4.48% | 0.48/12.03% | 0.25/23.81% | 3.75/47.17% | 0.00/0.74% | 0.17/7.32% |
| MBP(GID:4155)/**MBP | 7247 | | 0.07% | 0.08% | 0.01% | 0.01% | 0.02% | 0.03% | 0.25% | 1.40% | 0.10% | 0.78% | 0.00% | 0.03% |
| MMP1 | 5502/5502 | 0 | .40/37.99% | 0.04/0.89% | 0.49/11.03% | 0.08/12.47% | 0.02/1.40% | 0.06/6.17% | 0.19/7.71% | 0.10/2.53% | 0.29/27.42% | <u> </u> | 0.03/6.96% | 0.08/3.38% |
| CASP1 | 4554 | | 0.10% | 0.25% | 0.50% | 0.10% | 0.04% | 0.04% | <u>1.17%</u> | <u>13.24%</u> | 0.07% | <u>1.70%</u> | 0.01% | 0.03% |
| TNC(GID:3371)/**TNC | 3712 | | 0.01% | 0.02% | 0.31% | 0.01% | 0.03% | 0.00% | 0.02% | 0.11% | 0.08% | 0.17% | - | 0.05% |
| CTGF | 3607 | | 13.60% | 2.41% | 3.14% | 20.09% | 0.09% | 5.29% | 5.76% | 7.89% | 2.07% | 8.80% | 10.61% | 4.45% |
| CXCL1 | 3018 | | 0.17% | 0.12% | 0.63% | 0.14% | <u>11.37%</u> | 0.04% | 0.46% | 45.35% | 4.39% | 0.47% | 0.02% | 0.06% |
| SPARC | 2596 | | 0.02% | 0.01% | 0.56% | 0.00% | - | 0.03% | 0.01% | 0.02% | 0.01% | 0.06% | 0.01% | 0.00% |
| EPHX1 | 2326 | | 2.92% | 0.16% | 0.14% | 0.00% | 0.00% | 0.03% | 0.05% | 0.04% | 5.69% | 1.09% | 93.23% | 0.51% |
| 1.11 | 2004 | | 3.81% | 0.27% | 2.68% | 3.20% | 22.11% | 4.23% | 2.12% | 46.01% | 0.93% | 15.84% | 19.10% | 1.19% |
| FST | 1812 | | 1.21% | | 1.98% | 1.87% | 1.53% | 1.17% | 0.65% | 7.70% | | | 21.13% | 0.05 |

Summary

The Literature Lab[™] basic platform allows exploration of the PubMed database, highlighting co-occurences such as disease vs pathway, or co-morbidities *i.e.* disease vs disease. Breast, lung, and bone neoplasms were compared with pathways ranked according to strength of co-occurrence or LPF. The Literature Lab[™] PLUS platform allows the identification of statistically significant associations between gene lists and term domains in the PubMed literature. In the 70-gene signature from the van't Veer study, strong, moderate, and positive associations were linked with several hallmark pathways of metastasis, a result shared with the lung and bone gene signatures.

Interestingly, the gene set comparison feature highlighted several pathways specific to either lung or bone, or shared by both. This indicates that there are pathways specifically associated with these tumor metastatic microenvironments that may suitable for targeted investigations and potential therapeutic development. Moreover, there were several other pathway associations not typically associated with metastasis that may be fruitful for future investigations. This latter result highlights the hypothesis generation potential of the Literature LabTM platform.



The literature record within PubMed is too vast (>20 M publications to date) to permit comprehensive interrogation and identification of actionable associations. Over 4 M abstracts mention one or more human genes and modern high content genomic technologies are producing data at rates that outpace meaningful interpretation. Acumenta Biotech has created Literature LabTM, the only literature mining-based platform that identifies statistically significant associations between gene lists and key concepts in the literature. At the basic level, Literature LabTM can explore co-occurrences between term domains, *e.g.* diseases versus pathways. At a more rigorous level, Literature LabTM PLUS interrogates gene lists, including those derived via high content platforms, and scores the strength of each gene set / term domain association and significance based on 1000 random gene list comparison. It respects the uniqueness of each gene set and returns consistent and unique results. Literature LabTM identifies significant associations in a time efficient manner and reveals concepts and relationships in the literature that other gene analysis platforms cannot.

Damon Anderson, PhD Application Scientist 412-901-7785 danderson@acumenta.com The best way to see the power and benefits of Literature Lab[™] directly is to send us a gene list for complimentary analysis. <u>Click here</u> to register for a complimentary analysis or to be notified about upcoming Literature Lab Webinars.