

PET54 literature profiles

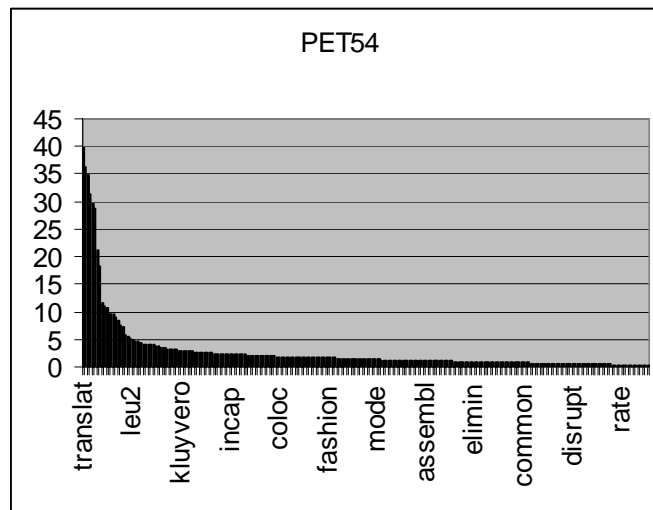
Literature relevant to a given gene might discuss several functional roles of the gene in the cell. To illustrate the differences in terms of literature profiles we compare the representation of genes using the vector space model, the clustered space model, and our own method (semantic profile). We provide the representations of a particular gene (*PET54*) from the SGD8 dataset analyzed in the manuscript, as well as the stemmed term sets that can be used to interpret each representation.

PET54 encodes a protein, located in the mitochondrial inner membrane, required for splicing the COX1 intron AI5 beta and translation of the COX3 mRNA.

1. PET54 in vector space representation (ordered by decreasing weights).

This representation is obtained from the gene-term frequency matrix calculated as described in the Methods section (i.e. it is the same gene-term frequency matrix used in NMF analysis)

Dimension: 2365. Only nonzero weight terms are shown.



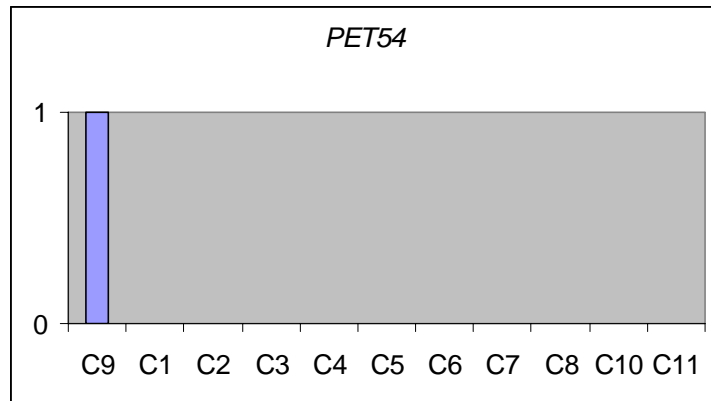
Gene profiles in the vector space model can be analyzed by the examination of most relevant terms (ordered by decreasing weights). E.g. top10 terms of *PET54* reveal important hints of the biological characterization of this gene.

translat
oxidas
mitochondri
cytochrom
leader
intron
mrna
beta
excis
splice

The vector space representation of this gene effectively accounts for several of its characteristics, as shown by the stemmed terms with highest weights (*translat*, *oxidas*, *mitochondri*, *cytochrom*, *leader*, *intron*, *mrna*, *beta*, *excis*, *splice*). However, the vector used for this representation has a very large dimension (dim=2,365).

2. *PET54* in clustered-based representation (ordered by decreasing weights).

In order to obtain an equivalent clustered space model using the vector space representation, we clustered the SGD8 dataset using k-means algorithm (employing the same gene-term frequency matrix used in NMF analysis). In a similar evaluation to that performed on the semantic features, we provided the top 10 terms of the centroids (ordered by decreasing weighting) of each of the 11 clusters to the four experts for interpretation. They were able to understand 8 of the 11 clusters. From the 3 remaining, only one of the experts labeled 2 clusters, while one cluster was not interpretable by none any of the experts.



Gene profiles in this space can be analyzed by the examination of most relevant terms of the corresponding cluster. E.g. the *PET54* profile can be interpreted by the top10 terms of the centroid in cluster #9.

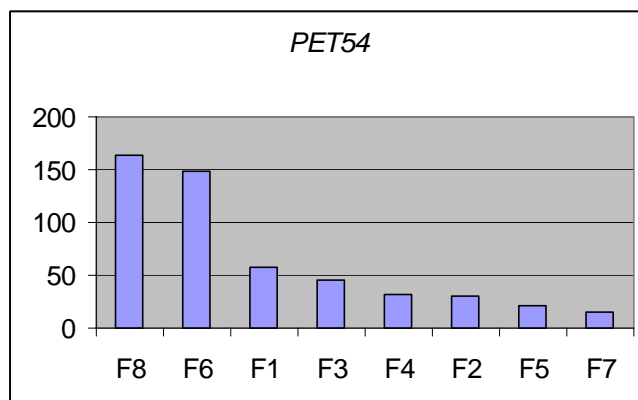
Table: Top10 stemmed terms corresponding to the centroids of the 11 clusters obtained by k-means algorithm. Labels shown if term set was meaningful to at least two of the four experts. Relevant feature for *PET54* is highlighted in red.

Clus1 <i>No label</i>	Clus2 <i>Chromatin remodeling</i>	Clus3 <i>Transport</i>	Clus4 <i>Lipid metabolism</i>	Clus5 <i>Mitosis</i>	Clus6 <i>No label</i>
synthetas fatti trna peroxisom acet ligas coli oxygen enzym ester	nucleosom swi histon snf chromatin remodel tbp acetyl h2a rsc	v-atpas vacuolar golgi vesicl vacuol copper transport membran iron snare	sterol sphingolipid ergosterol inositol ceramid phospholipid reductas phosphatidylcholin phosphatidylserin cytochrom	spindl kinetochor cyclin centromer checkpoint microtubul anaphas mitosi spb mitot	ribosom actin translat pheromon gcn4 alpha mate kinas wall calcineurin

Clus7 <i>No label</i>	Clus8 <i>Transport</i>	Clus9 <i>Mitochondria</i>	Clus10 <i>DNA replication</i>	Clus11 <i>DNA repair</i>
glucos permeas transport hexos nitrogen uptak camp snf1 trehalos ubiquitin	autophagi copii vesicl autophagosom conjug coat cargo vacuol membran golgi	mitochondri inner preprotein export mrna mitochondria transloc hsp70 outer membran	replic dna telomer repair orc dsb pcna rfc checkpoint mcm	repair ner tfiih excis dna mismatch rad6 damag apn1 rad1

3. PET54 in NMF representation (ordered by decreasing weights).

Dimension: 8.



NMF representation was obtained as described in the manuscript.

The representation of *PET54* by means of NMF involves the linear combination of 8 features. Gene profiles in this space can be analyzed by the examination of most relevant terms in each feature. This profile shows highly significant weights to both the *mitochondria* and the *protein synthesis* features, revealing a more complete description of the gene in terms of local features. Although this particular gene was finally assigned to cluster *G (mitochondria)*, similar to the cluster created by k-means, its literature profile provides a more detailed representation using a reduced set of features.

Table: Top 10 stemmed terms in the k=8 semantic features obtained for a NMF experiment (ordered by decreasing importance). Labels show topical interpretations provided by experts (including more concrete topics in parenthesis). Relevant features for PET54 are highlighted in red

F1 DNA metabolism (DNA replication)	F2 DNA metabolism (DNA repair)	F3 Metabolism / Stress / Degradation	F4 Transcription (chromatin)	F5 Cell Division (mitosis)	F6 Mitochondria	F7 Transport (vesicular trafficking)	F8 Protein synthesis
Replic pcna dna ner damag checkpoint rfc pol polymeras rad6	repair telomer dsb recombin mismatch dna rad52 excis rad51 endonucleas	glucos fatti heat stress endoplasm reticulum proteasom phosphatas atpas sphingolipid	actin swi nucleosom snf histon chromatin elong mate silenc polar	spindl cyclin kinetochor hsp90 chaperon scf anaphas mitosi centromer mitot	mitochondri preprotein mitochondria inner transloc outer membran matrix oxid translocas	transport vesicl vacuolar vacuol membran nitrogen secretori autophagi cytoplasm sort	translat mrna trna alpha gcn4 beta gtp phosphoryl exchang kinas