**Additional file 2: Supplementary figures for**

**Comparative genomics and transcriptomics of *Pichia pastoris***
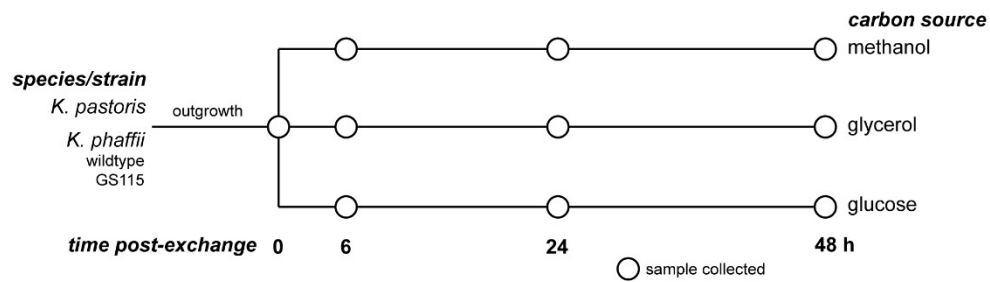


**Figure S1: Experimental timeline for RNA sequencing of *Komagatella* strains.**

Schematic timeline for collection of RNA samples during batch cultivation of strains in shake flasks on 3 different carbon sources.  Three independent cultivations were sampled for each time point.
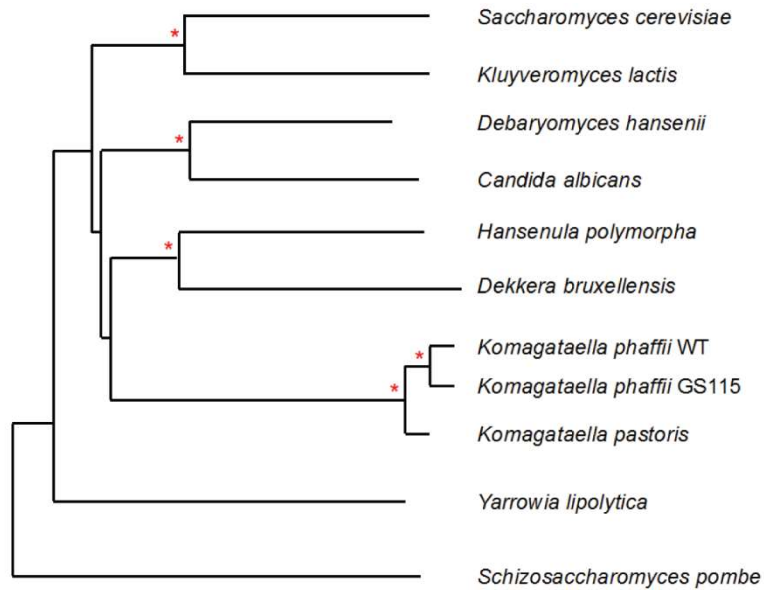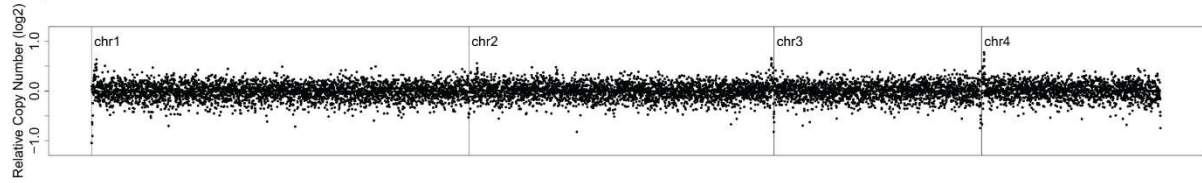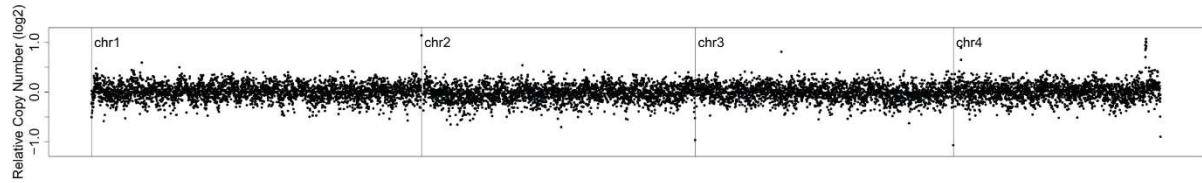
**Figure S2: Phylogenetic comparison of *K. pastoris* and *K. phaffii* to other related yeasts.** Phylogeny was generated using a concatenated, gap-free alignment of ten orthologous proteins. Phylogenetic tree was calculated using neighbor-joining, distance-based, maximum likelihood and maximum parsimony methods; reliability was assessed using bootstrapping. Clades marked with an asterisk are supported by 100% of bootstrap replicates in all four methods.

**Figure S3: Gene conservation between *K. pastoris* and *K. phaffii*.** a) Histogram denoting homology at the base pair level for all 1:1 orthologous genes. b) Alignment of the P$_{GAPDH}$ promoter element between *K. pastoris* and *K. phaffii*. c) Alignment of the P$_{AOX1}$ promoter element between *K. pastoris* and *K. phaffii*.

a) *K. pastoris*
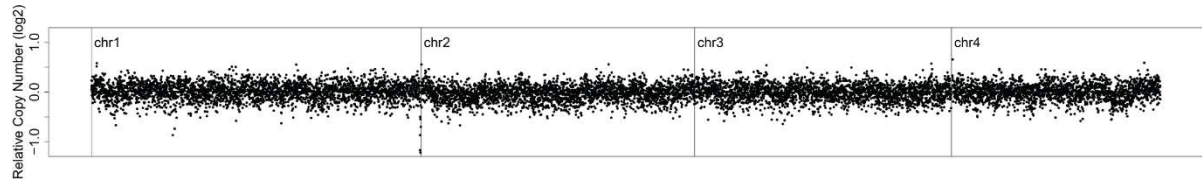


b) *K. phaffii*



c) *K. phaffii (GS115)*



**Figure S4:** Copy number determination for major chromosomes in a) *K. pastoris*, b) *K. phaffii* wild-type and c) *K. phaffii* GS115 strains.

a) *K. pastoris*



b) *K. phaffii*



**Figure S5:** Codon usage for a) *K. pastoris* and b) *K. phaffii* as determined from all coding sequences identified in genome annotation. The relative abundance observed for each codon is represented as a percentage of total codon usage for the corresponding amino acid.

**Figure S6: Isoform expression in *K. pastoris* and *K. phaffii* as a function of cultivation conditions.** Heat maps of gene expression (log2 fpkm) for isoforms of alternatively spliced genes that alter coding sequences in a) *K. pastoris* and b) *K. phaffii* detected in initial genome annotation. Alternatively spliced genes with sufficient homology to *S. cerevisiae* are named, otherwise gene identifiers from genome annotation are used. Isoform expression is shown as a function of batch growth in glycerol, glucose or methanol during a 48 h cultivation period.

**Figure S7: Chromosomal locations of highly expressed genes.** Map of chromosomal location (base pair identity) for the most highly expressed genes (top 10% expression) in a) *K. pastoris* and b) *K. phaffii*. Black lines indicate gene expression level at 24 h time points during batch cultivation in either glycerol, glucose or methanol. Red lines indicate locations of GC-rich autonomously replicating sequence (GC-ARS) motifs identified by BLAST.
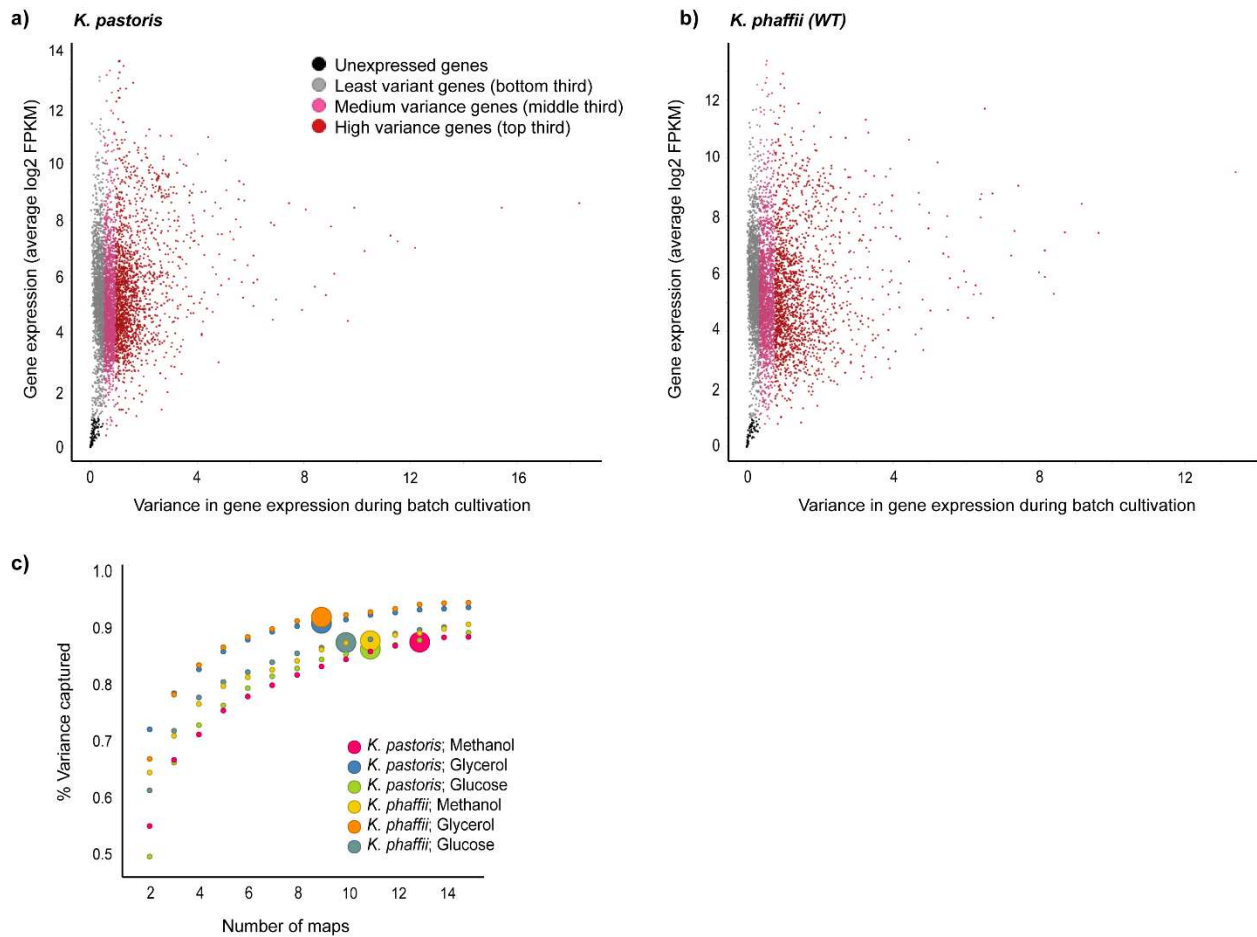
**Figure S8: Variance in gene expression during batch cultivation of *K. pastoris* and *K. phaffii*.** Scatter plots of the average variance versus expression observed for all annotated genes across the 10 conditional averages generated from either a) *K. pastoris* or b) *K. phaffii* expression data. Genes with average log2 fpkm <1 and variance < 0.05 were excluded from further analyses. c) Elbow analysis of input cluster number to identify optimal expression data clustering by self-organizing maps (SOMs). Large circles denote the number of clusters for each expression data set where the additional variance captured by further clustering was < 1%.

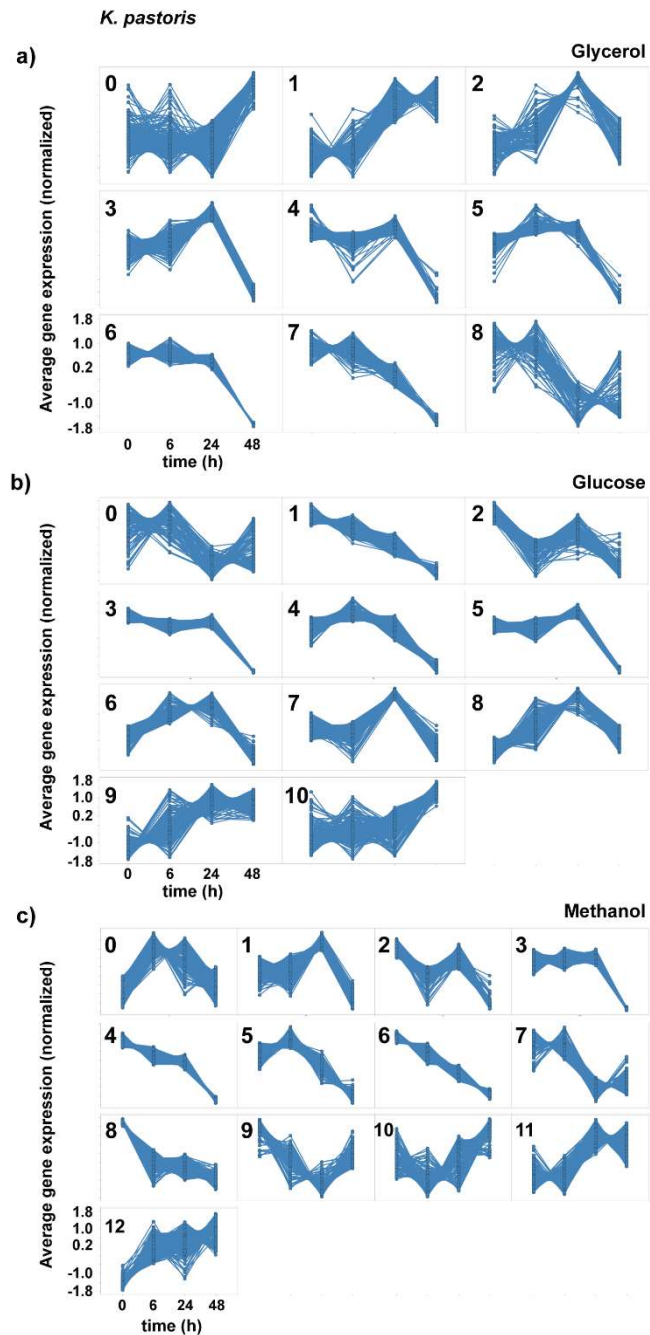**Figure S9: Gene expression phenotypes in *K. pastoris* as a function of cultivation conditions.** Self-organizing maps (SOMs) of genes changing expression similarly in *K. pastoris* during a 48 h batch cultivation in a) glycerol, b) glucose or c) methanol.
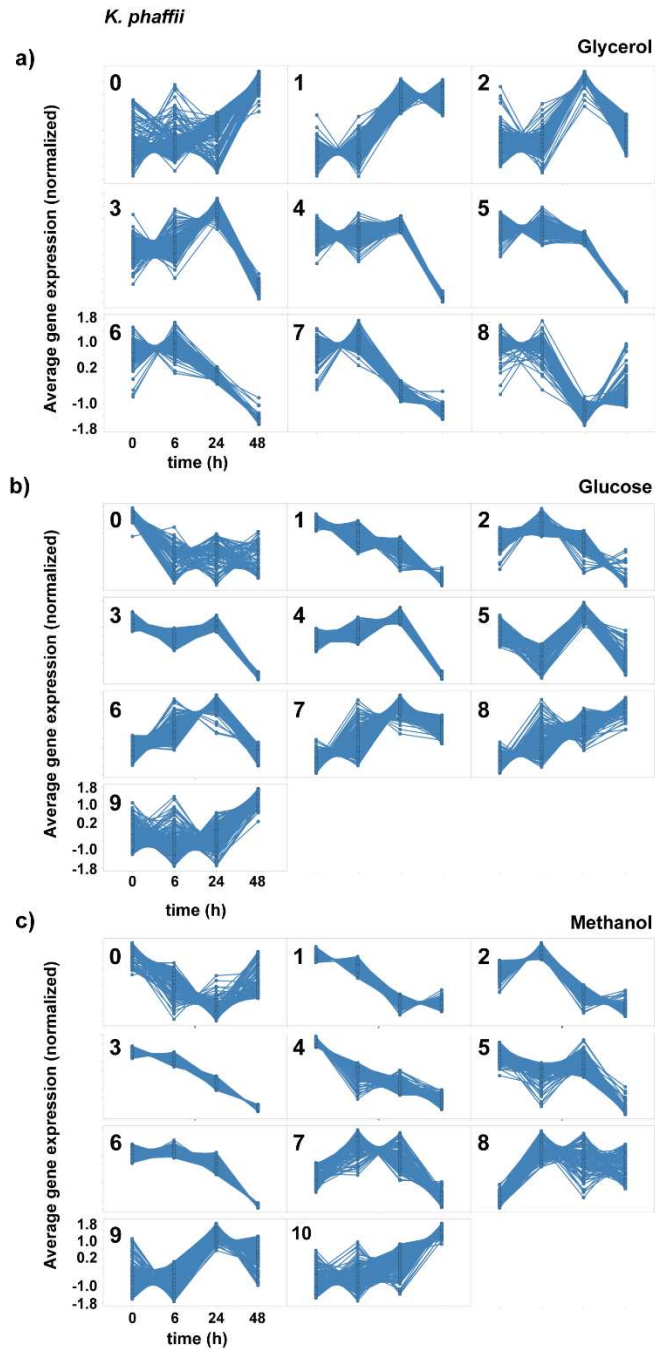
**Figure S10: Gene expression phenotypes in *K. phaffii* as a function of cultivation conditions.** Self-organizing maps (SOMs) of genes changing expression similarly in *K. phaffii* during a 48 h batch cultivation in a) glycerol, b) glucose or c) methanol.
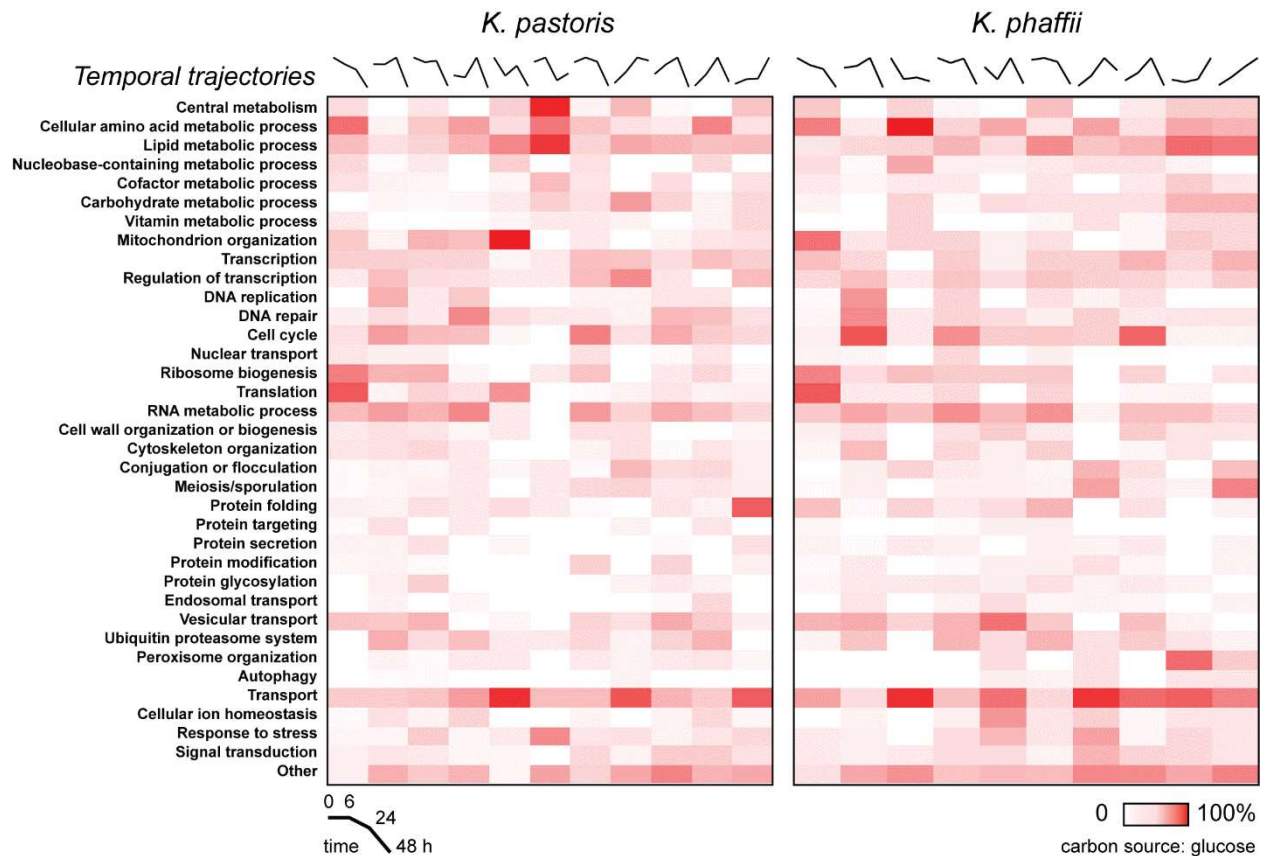
**Figure S11: Biological process enrichment as a function of cultivation in glucose.** Heat map representation of the enrichment of GO biological process terms for expression phenotypes observed in *K. pastoris* and *K. phaffii* during a 48 h batch cultivation in glucose as characterized by self-organizing maps (SOMs). Representative temporal trajectories of gene expression were generated for each SOM by averaging expression data at each time point for genes present within a given map. Color density relates to the number of genes assigned to a particular process as a percentage of the total number genes present in a particular expression phenotype or map.

**Figure S12: Biological process enrichment as a function of cultivation in glycerol.** Heat map representation of the enrichment of GO biological process terms for expression phenotypes observed in *K. pastoris* and *K. phaffii* during a 48 h batch cultivation in glycerol as characterized by self-organizing maps (SOMs). Representative temporal trajectories of gene expression were generated for each SOM by averaging expression data at each time point for genes present within a given map. Color density relates to the number of genes assigned to a particular process as a percentage of the total number genes present in a particular expression phenotype or map.
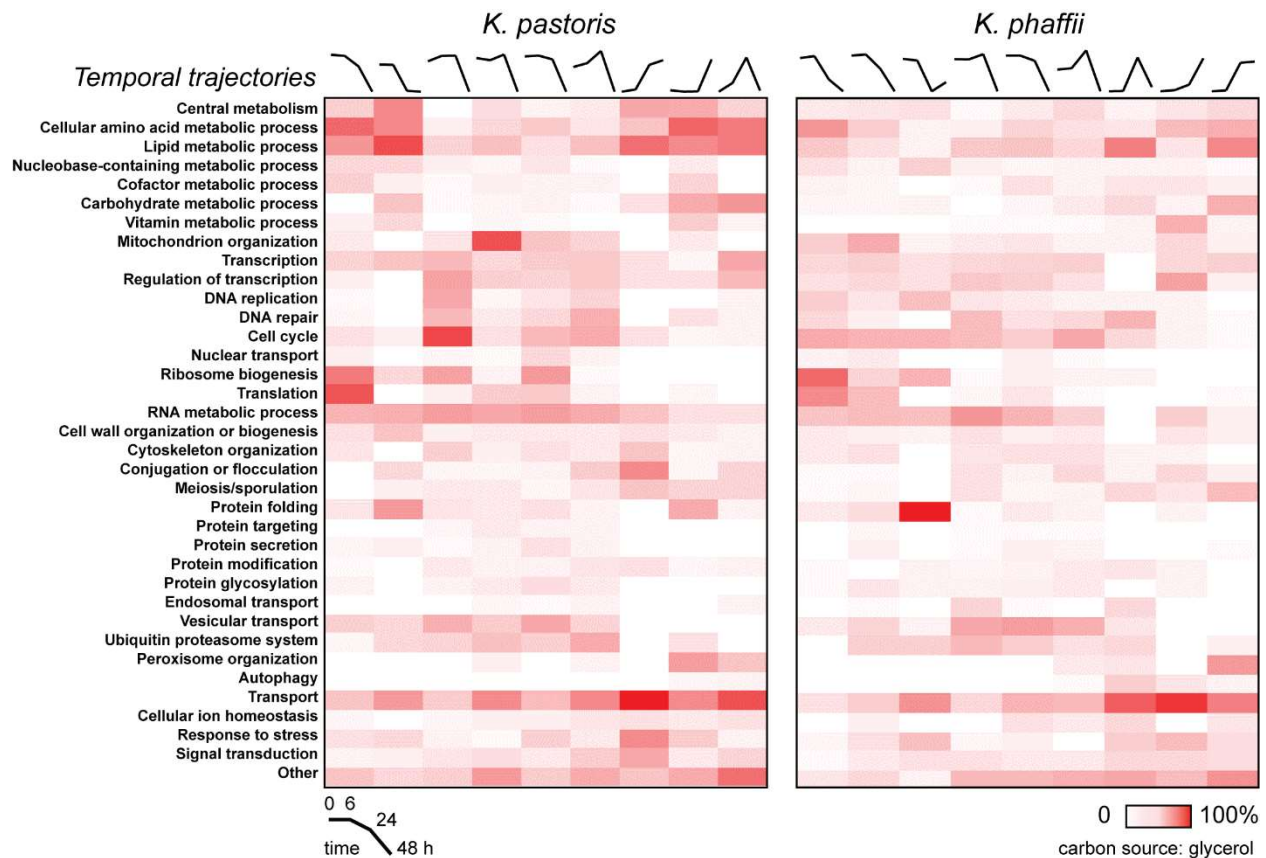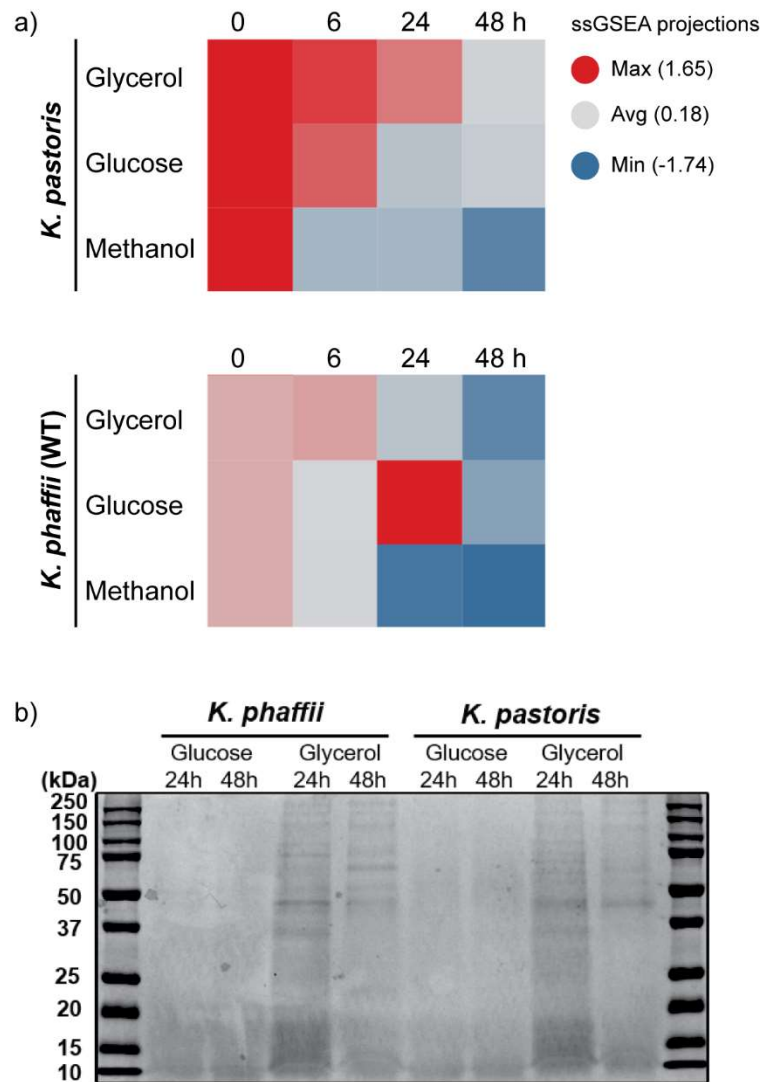
**Figure S13:  Secretory pathway protein expression in _K. pastoris_ and _K. phaffii_.**

a) Row normalized single set Gene Set Enrichment Analysis (ssGSEA) projections for 170 proteins bearing a signal peptide as identified by Signalp.  b) SDS-PAGE analysis of host-cell protein expression in supernantants during batch cultivation of _K. pastoris_ and _K. phaffii_ for 48h in glucose or glycerol-containing media.
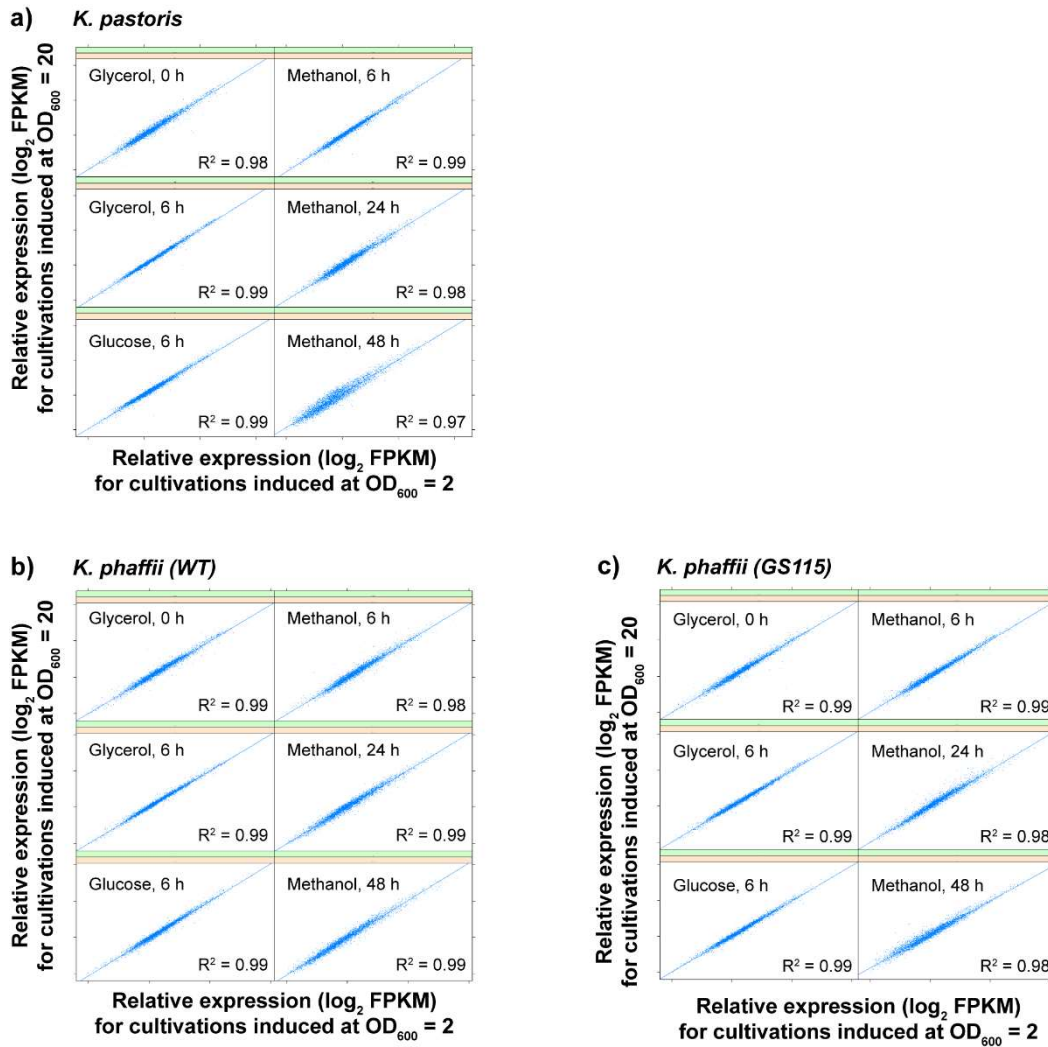
**Figure S14: Correlation of gene expression at two different cultivation densities.**

Scatter plots of gene expression between similar cultivation conditions for a) *K. pastoris*,

b) wildtype *K. phaffii*, and c) *K. phaffii* GS115 grown at two different cell densities.

Density A corresponds to cultures outgrown to $OD_{600}$ = 2.0 prior to sampling and

Density B corresponds to cultures outgrown to $OD_{600}$ = 20 prior to sampling. Pearson

correlation coefficients were calculated from expression vectors that were averages of

three biological replicates for each cultivation condition and density.