

SUPPORTING INFORMATION

Table S1 Expert pathologist annotations collected for model training by substance.

Annotation numbers refer to pixel level bounding boxes (regions) drawn by pathologists using the PathAI digital slide platform. Bolded annotations reflect those used as classes during model training.

Table S2 Manual scoring data for NAS components and fibrosis staging supporting intra-reader and inter- reader variability. Slide IDs indicate unique patient biopsies.

Rows with repeated IDs were used to assess intra-pathologist reproducibility. Reading intervals specify time between reads of H&E and Trichrome slides.

Table S3 Intra-observer reproducibility of pathologist manual scores of NAS components and fibrosis staging. Values are weighted Cohen's kappa (1 is perfect agreement) for each pathologist's repeated measurements. Cohen's Kappa measures the rate of agreement while accounting for the degree of disagreement based on the ordinal gradation.

Table S4 Inter-observer reproducibility of fibrosis staging. Values are weighted Cohen's Kappa for each pathologist's and the model's score against the pathologist consensus. For each pathologist their score was computed against the consensus of the other

pathologists.

Table S5 Prognostic value of each ML feature for predicting progression to cirrhosis and adjudicated clinical events in STELLAR-3 and STELLAR-4, respectively. Each row represents results from a univariate Cox regression model. Confidence intervals (ci-low, and ci-high) represent the 95% confidence interval for the hazard ratio.

Table S6 Prognostic significance of feature clusters for prediction of adjudicated clinical events and progression to cirrhosis. Feature column lists all features assigned to each cluster via agglomerative clustering. Q column lists the Benjamini-Hochberg corrected p-values. P column lists nominal cluster-wise p-values derived from EBM and univariate Cox regression.

Figure S1 Confusion matrix indicating model performance for labeling held-out pathologist annotations on H&E slides. Each column in the x-axis represents the predicted label and each column in the y-axis represents the true label. Box colors indicate the accuracy (0.0, white to purple, 1.0). Normal liver represents an amalgam of normal liver histologic features including normal hepatocytes (**Table S1**).

Figure S2 Associations between ML-based model measurements and grading of NAS features by the CP according to trial. Model values describe the percentage of tissue predicted to be the substance in question (steatosis, lobular inflammation, or hepatocellular ballooning). Boxes show the IQR and whiskers show 1.5 x the limit of the IQR. Points show values beyond this range. Values shown are Spearman correlations (ρ) and corresponding p-values. No models were trained using the ATLAS dataset.

Figure S3 Confusion matrix indicating model performance for labeling held-out pathologist annotations on trichrome-stained slides. Each column in the x-axis represents the predicted label and each column in the y-axis represents the true label. Box colors indicate the accuracy (0.0, white to purple, 1.0). Confusion between starred classes (Blood Vessels and Lumen) are expected since these annotations are overlapping.

Figure S4 Associations between ML-based model measurements and staging of fibrosis by the CP according to A) the Ishak classification in all studies, and B) the Ishak and NASH CRN classifications by trial. Model values describe the weighted average score for each patient by pathologist-derived fibrosis stage. The weighted average score is computed by multiplying each fraction of tissue area by its corresponding stage and summing across the slide (**Methods**). Boxes show the IQR

and whiskers show 1.5 x the limit of the IQR. Points show values beyond this range. Values shown are Spearman correlations (ρ) and corresponding p-values. No models were trained using the ATLAS dataset.

Figure S5 Reproducibility and variability of Ishak fibrosis staging by pathologists in 166 slides. A) Bar chart showing intra-rater agreement (weighted Cohen's kappa) for Ishak fibrosis stage. B) Bar chart showing inter-rater agreement (weighted Cohen's kappa) for Ishak fibrosis stage. Each pathologist's agreement is measured with all other pathologists in the study.

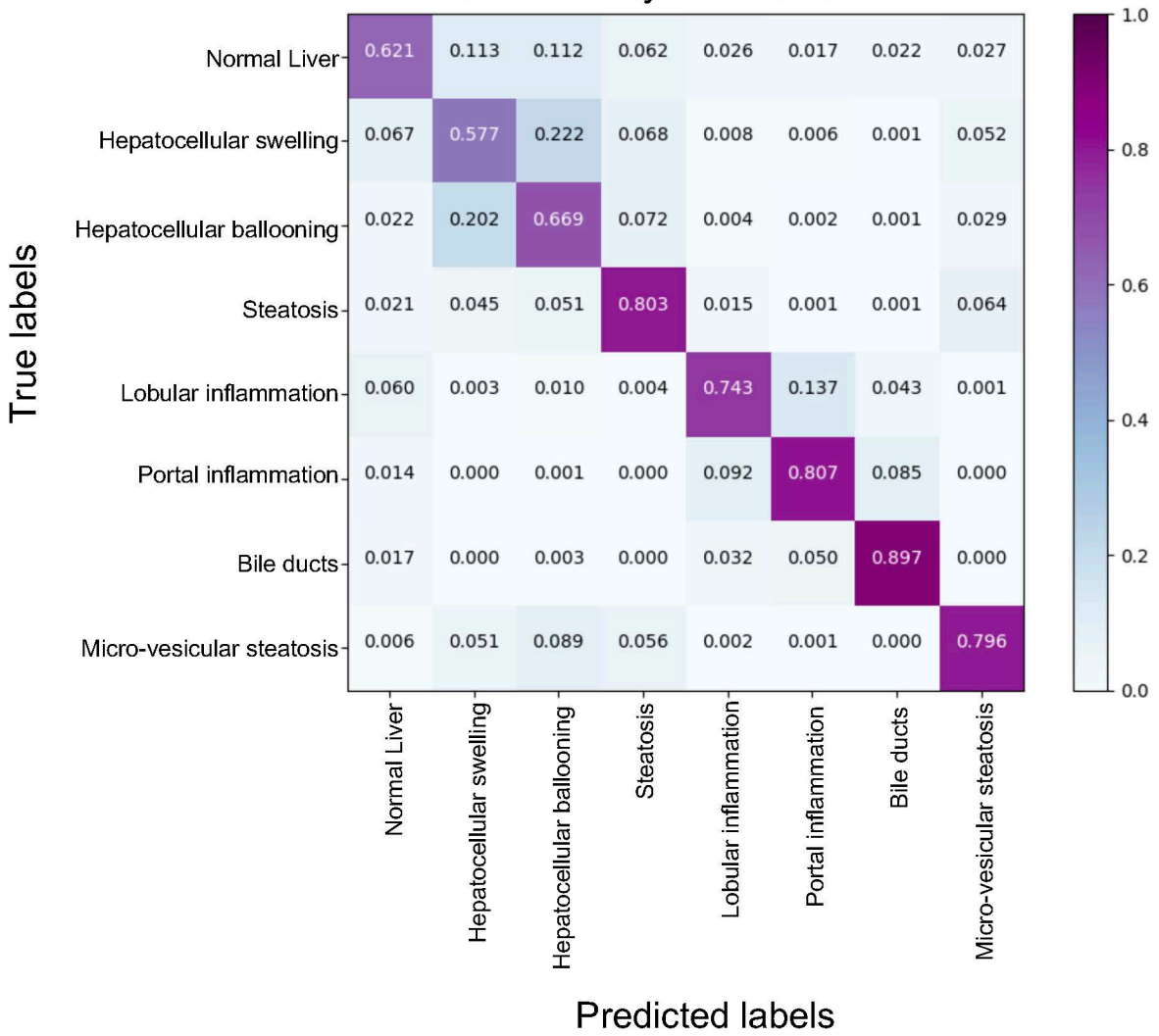
Figure S6 Feature clustering and significance testing for patient prognosis in STELLAR-3 and STELLAR-4. A) Heatmap showing the inter-feature absolute Spearman correlation (0=black, white=1). Features are sorted using agglomerative clustering. B) "Elbow" plot showing the number of clusters by height of the dendrogram in (A). Blue line indicates selected cut point corresponding to 22 clusters. C and D) Quantile-quantile plots of cluster-wise nominal p-values (blue dots) for association with progression to cirrhosis in patients with bridging (F3) fibrosis (STELLAR-3) or liver-related clinical events in patients with cirrhosis (F4; STELLAR-4). Significant clusters (FDR corrected q-value) are denoted by red circles. P-values are sorted and plotted against the corresponding uniform quantile value.

Figure S7 Correspondence of DELTA Liver Fibrosis score with markers of treatment response. A-C) Box and whisker plot showing the difference in DELTA Liver Fibrosis score for patients in the ATLAS trial by treatment group (x-axis) and achievement of a ≥ 1 -stage improvement in NASH CRN fibrosis stage as evaluated by the CP (A), reduction of at least 0.5 in ELF score (B) and 25% reduction in FibroScan (C). Boxes show the IQR and whiskers show 1.5 x the limit of the IQR. P-values are computed using the Mann-Whitney test.

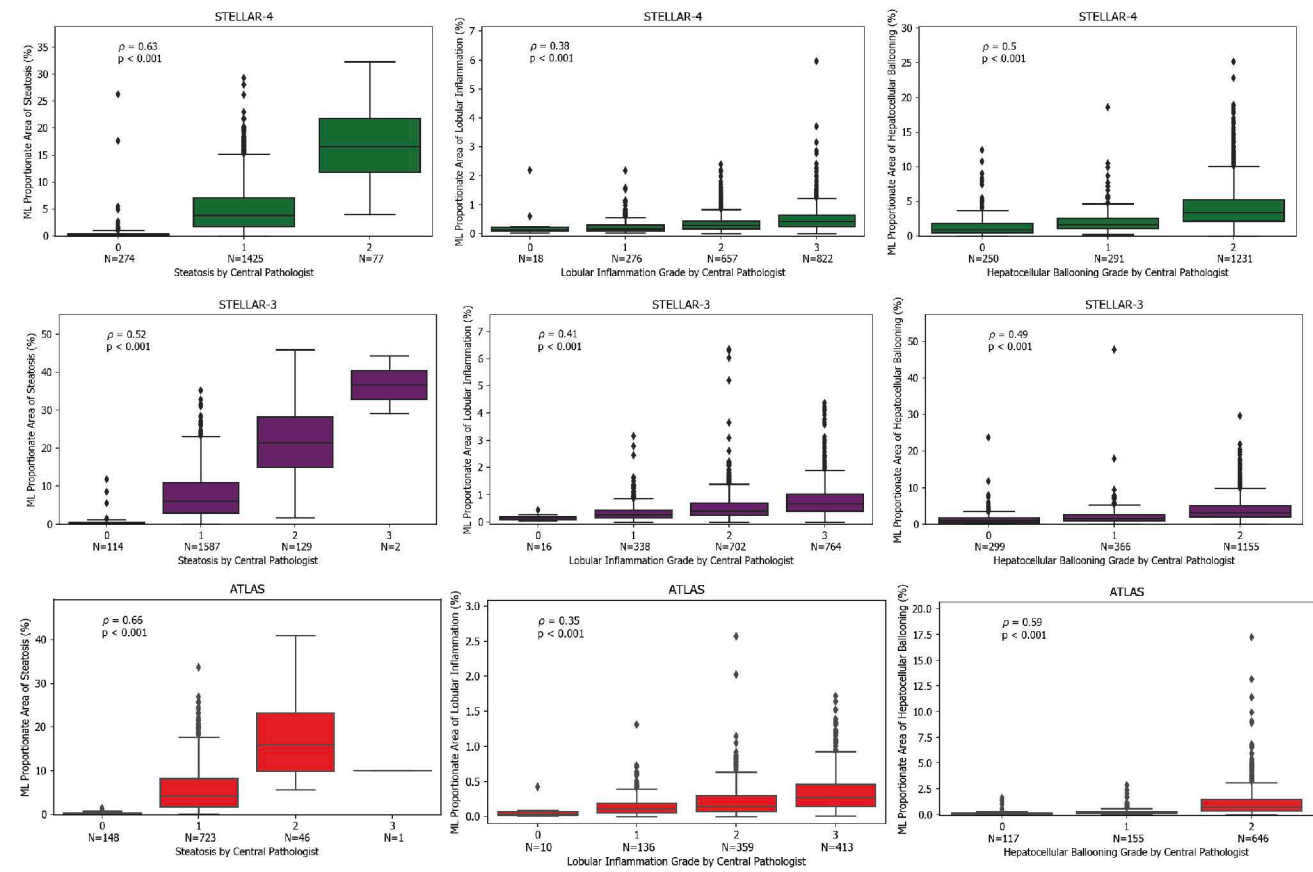
Figure S8 Bar charts comparing the proportion of patients with a reduction in fibrosis as assessed by the DELTA Liver Fibrosis score and the CP in experimental arms (red) compared with placebo-treated patients (grey) in ATLAS. P-values are computed using Fisher's exact test.

S. Figure 1

H&E Model accuracy on held out annotations



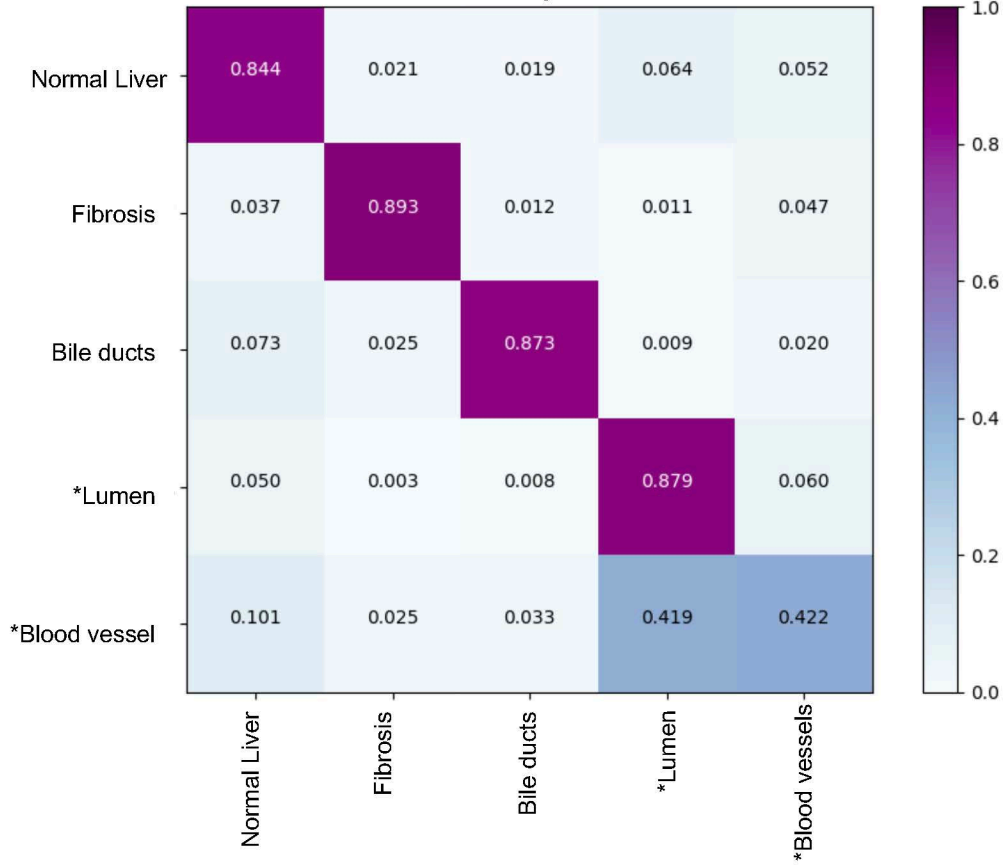
S. Figure 2



S. Figure 3

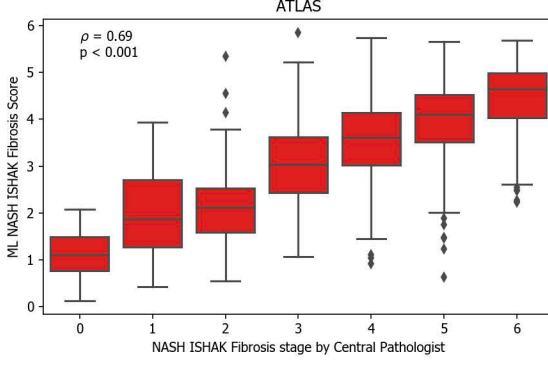
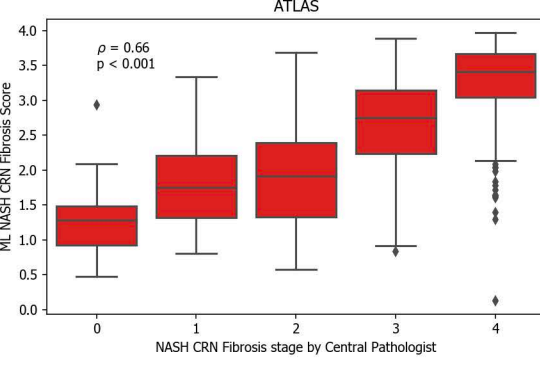
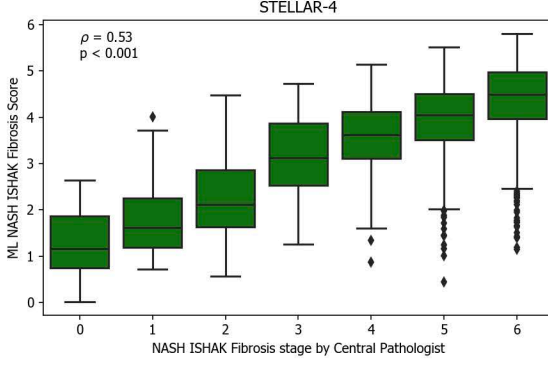
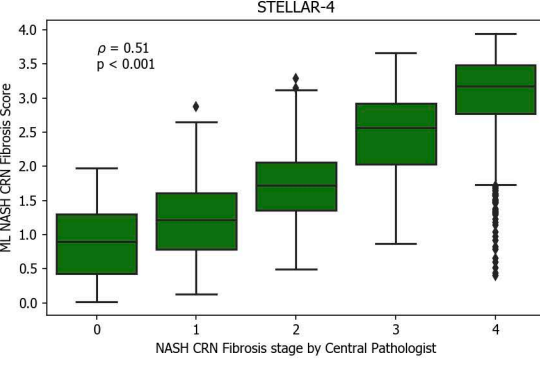
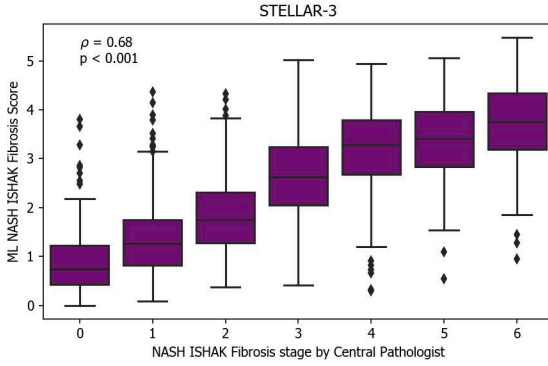
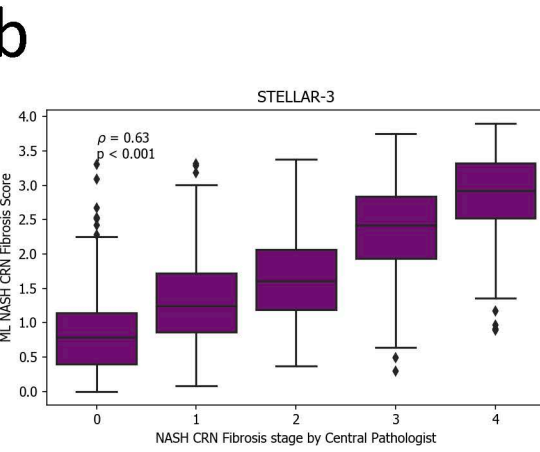
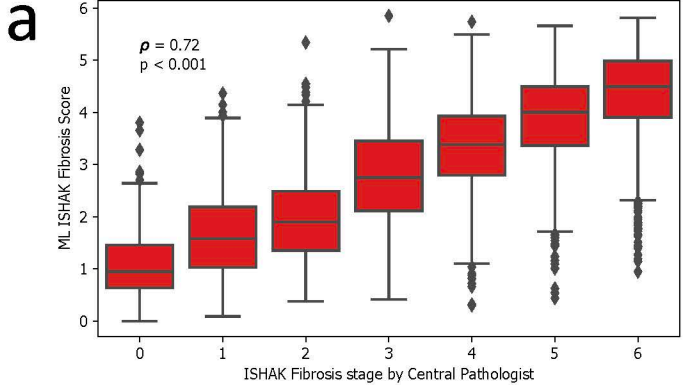
Trichrome Model accuracy on held out annotations

True labels



Predicted labels

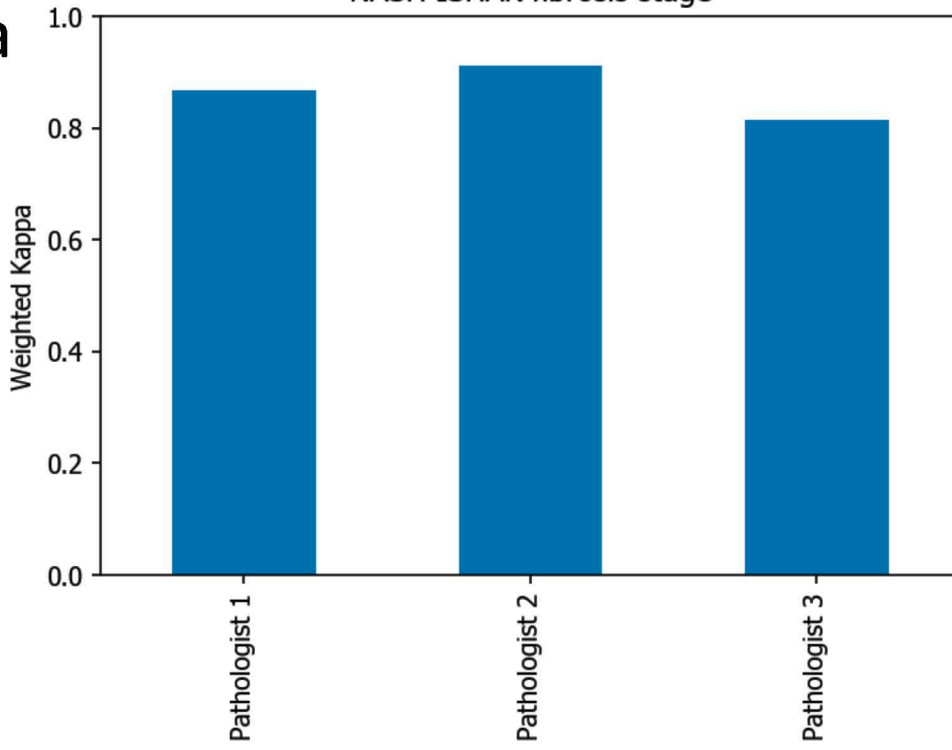
S. Figure 4



S. Figure 5

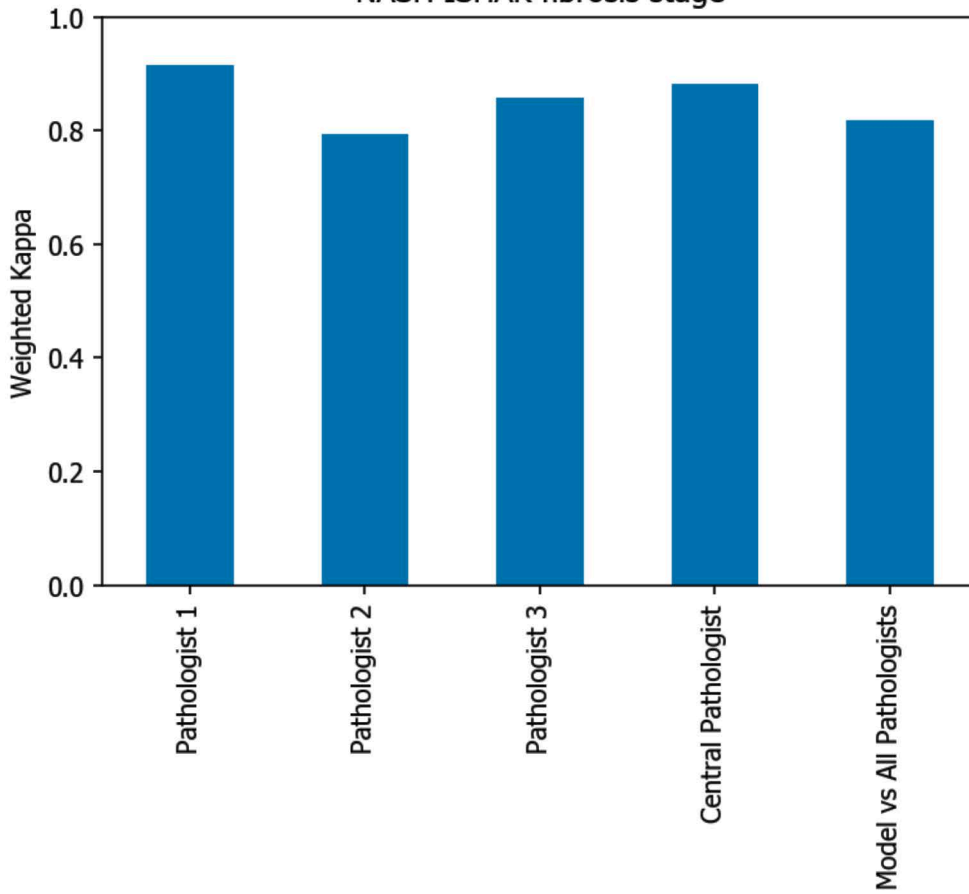
NASH ISHAK fibrosis stage

a

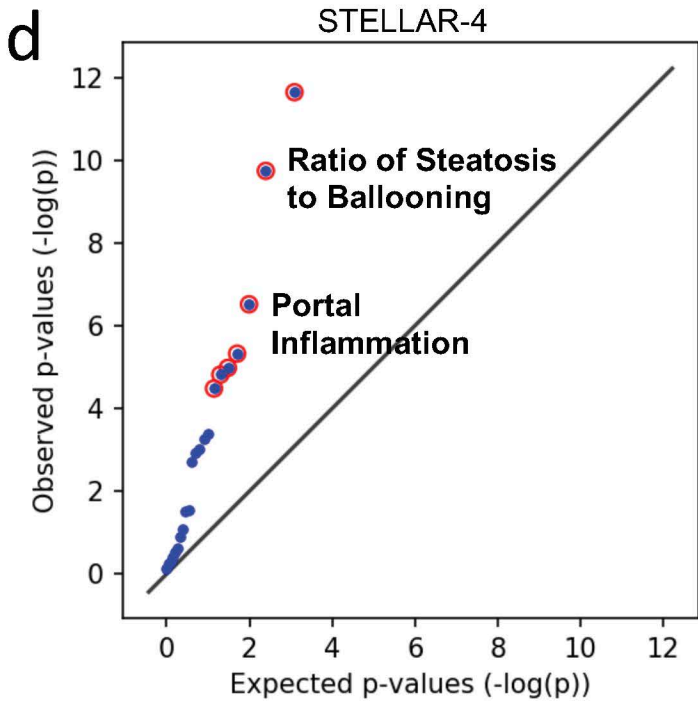
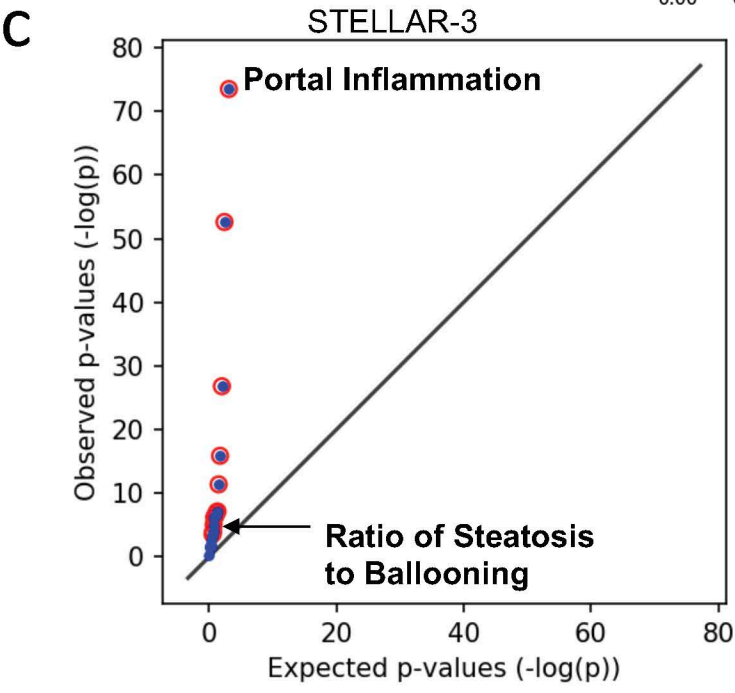
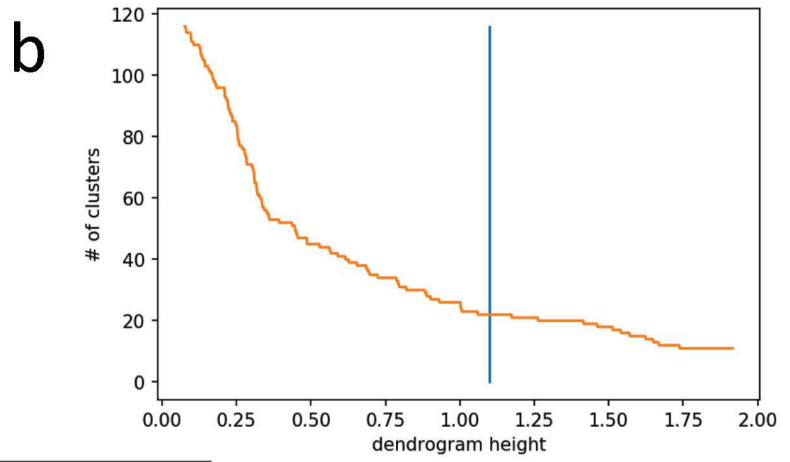
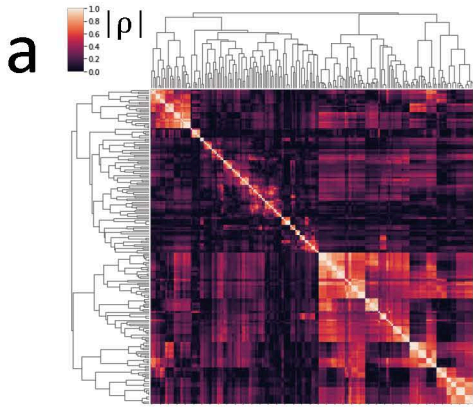


b

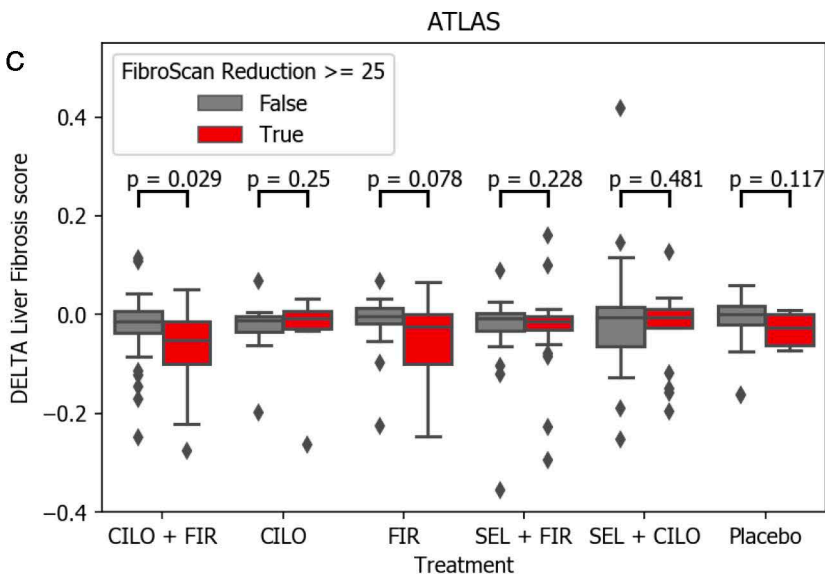
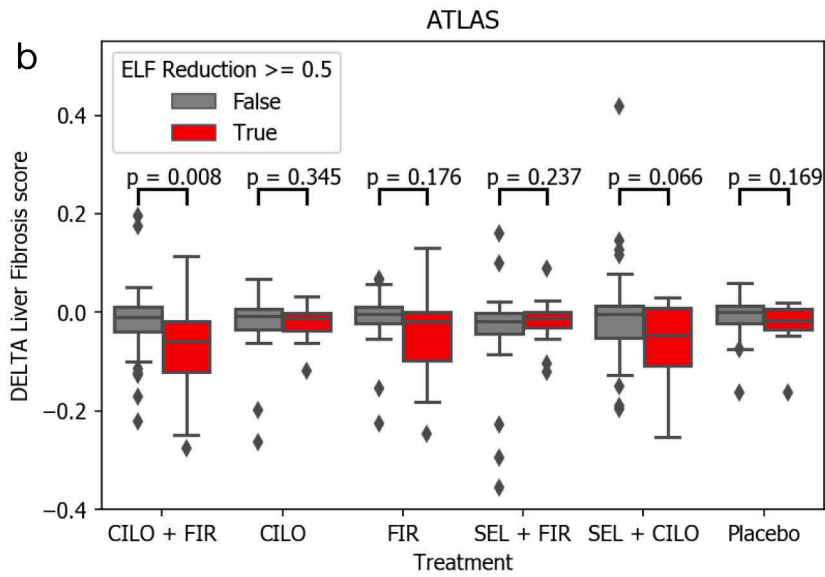
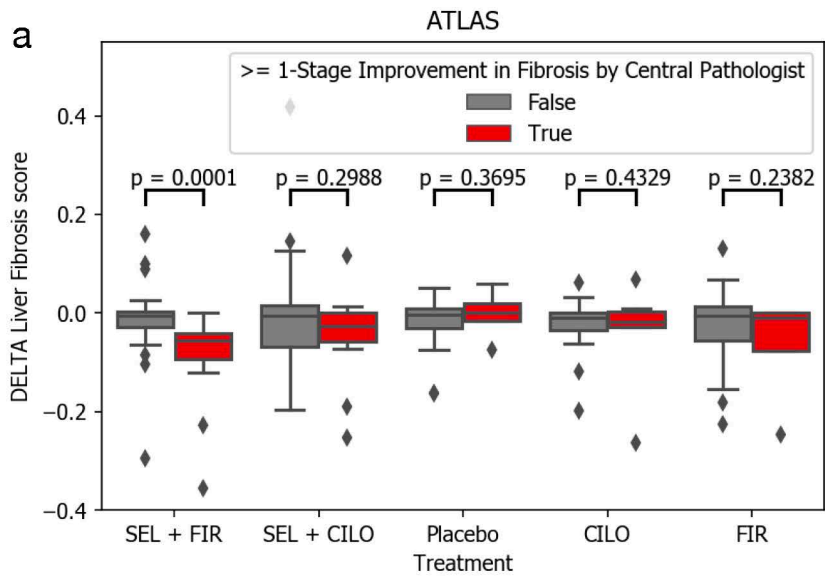
NASH ISHAK fibrosis stage



S Figure 6



S Figure 7



S Figure 8

