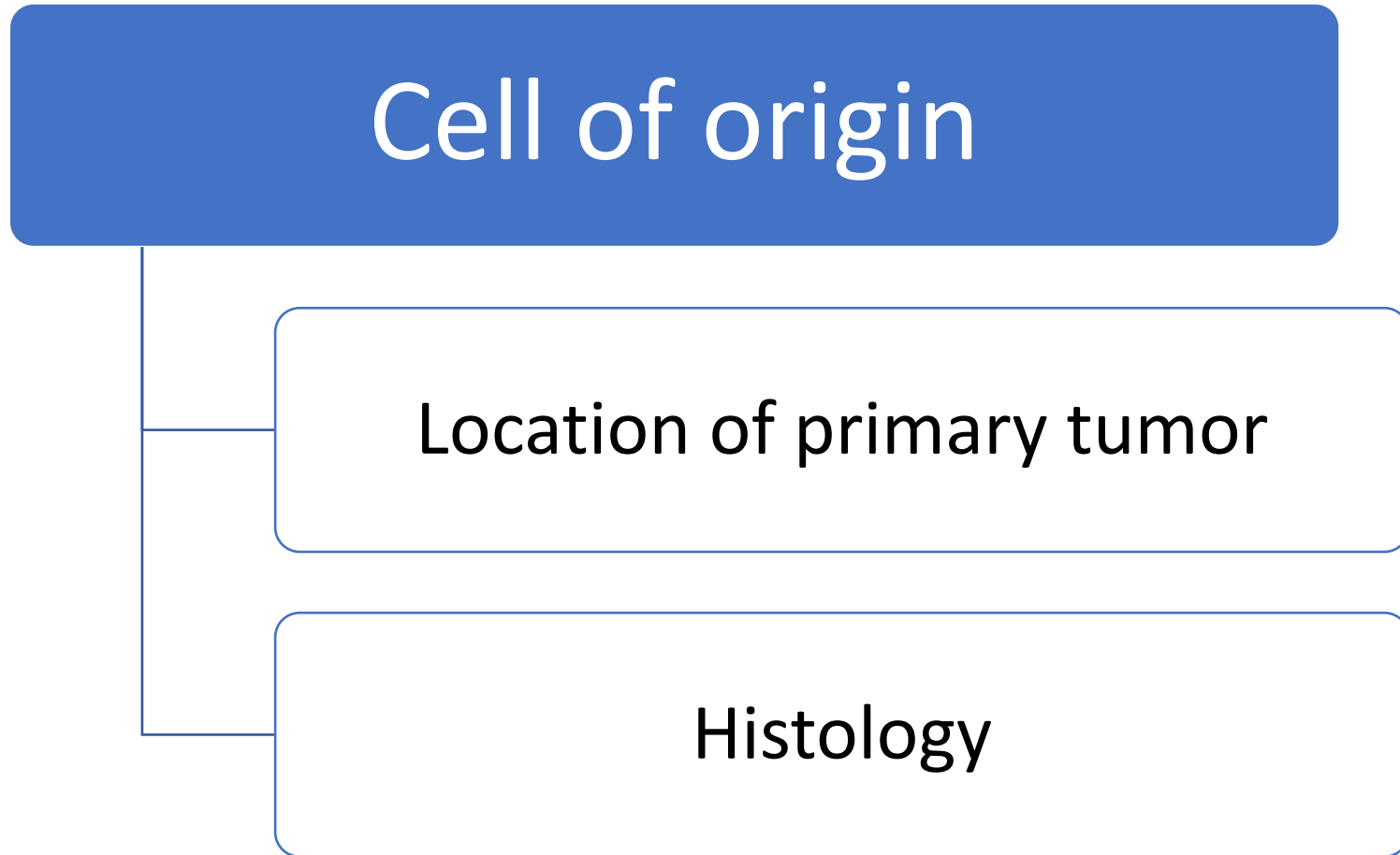# Identifying tumor's organ of origin using deep learning

Mohammed El-Kebir

# Major determinant of outcome

# Datasets

**Table 1 Distribution of tumour types in the PCAWG training and test data sets.**

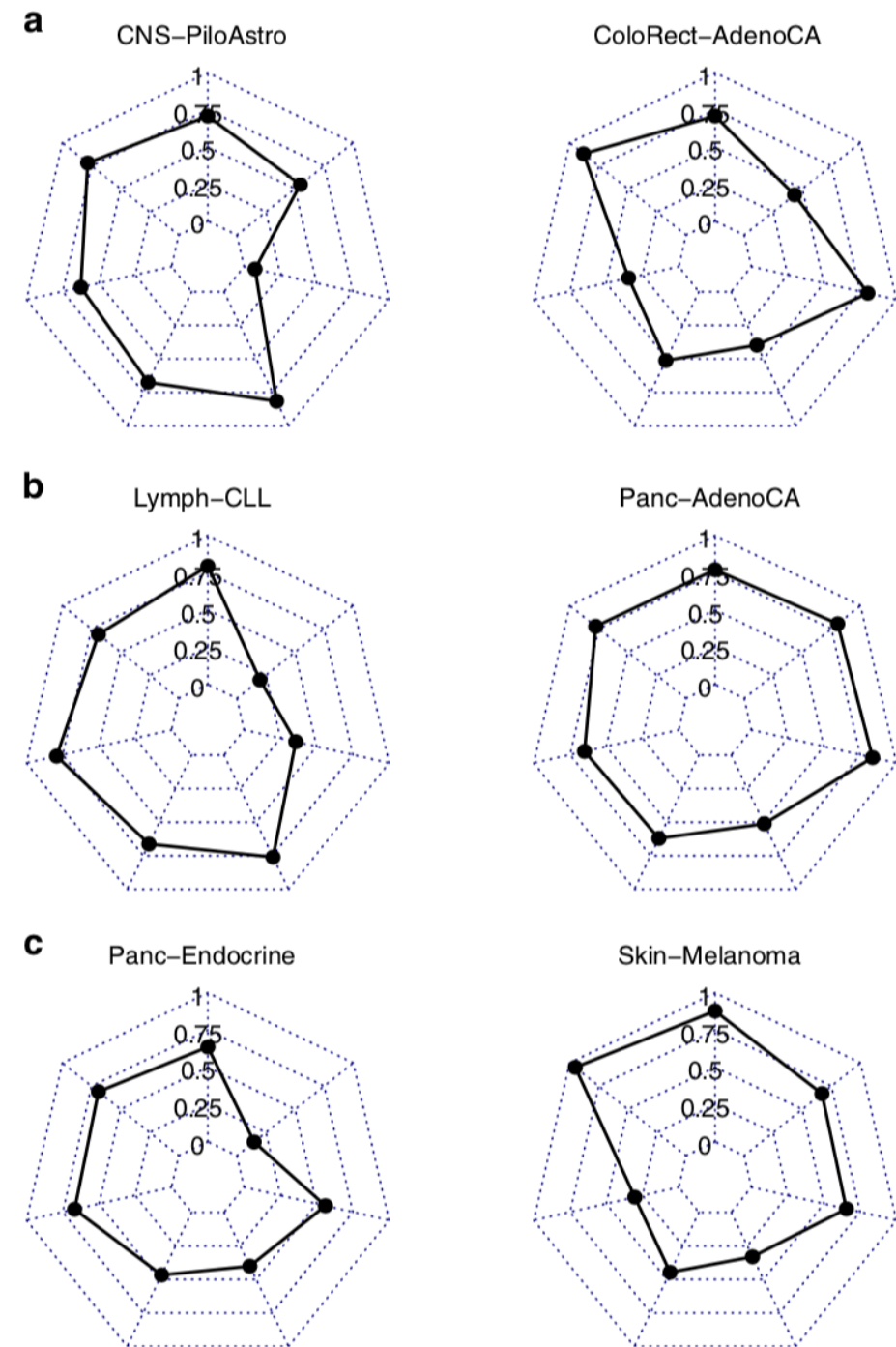| Abbreviation | Organ system | Tumour type | Tumour samples |
|---|---|---|---|
| Liver-HCC | Liver | Liver hepatocellular carcinoma | 306 |
| Panc-AdenoCA | Pancreas | Pancreatic adenocarcinoma | 235 |
| Breast-AdenoCA | Breast | Breast adenocarcinoma | 198 |
| Prost-AdenoCA | Prostate gland | Prostate adenocarcinoma | 189 |
| CNS-Medullo | Brain, cranial nerves and spinal cord | Medulloblastoma | 146 |
| Kidney-RCC | Kidney | Renal cell carcinoma (proximal tubules) | 143 |
| Ovary-AdenoCA | Ovary | Ovarian adenocarcinoma | 112 |
| Skin-Melanoma | Skin | Skin-melanoma | 106 |
| Lymph-BNHL | Lymph nodes | Mature B-cell lymphoma | 105 |
| Eso-AdenoCA | Oesophagus | Oesophageal adenocarcinoma | 98 |
| Lymph-CLL | Blood, bone marrow and hematopoietic sysstem | Chronic lymphocytic leukaemia | 95 |
| CNS-PiloAstro | Brain, cranial nerves and spinal cord | Pilocytic astrocytoma | 89 |
| Panc-Endocrine | Pancreas | Pancreatic neuroendocrine tumour | 85 |
| Stomach-AdenoCA | Stomach | Gastric adenocarcinoma | 70 |
| Head-SCC | Gum, floor of mouth and other mouth | Head/neck squamous cell carcinoma | 57 |
| ColoRect-AdenoCA | Large intestine (excluding appendix) | Colorectal adenocarcinoma | 52 |
| Lung-SCC | Lung and bronchus | Lung squamous cell carcinoma | 48 |
| Thy-AdenoCA | Thyroid gland | Thyroid adenocarcinoma | 48 |
| Myeloid-MPN | Blood, bone marrow and hematopoietic system | Myeloproliferative neoplasm | 46 |
| Kidney-ChRCC | Kidney | Renal cell carcinoma (distal tubules) | 45 |
| Bone-Osteosarc | Bones and joints | Sarcoma, bone | 44 |
| CNS-GBM | Brain, cranial nerves and spinal cord | Diffuse glioma | 41 |
| Uterus-AdenoCA | Uterus, nos | Uterine adenocarcinoma | 40 |
| Lung-AdenoCA | Lung and bronchus | Lung adenocarcinoma | 38 |
| | | | **2436** |

# Feature types

**Table 2 WGS feature types used in classifiers.**

| Feature category | Feature type | Feature count | Description |
|---|---|---|---|
| Mutation distribution | SNV-BIN | 2897 | Number of SNVs per 1-Mbp bin, and per chromosome, normalised against the total number of SNVs per sample |
| | CNA-BIN | 2826 | Number of CNAs per 1-Mbp bin |
| | SV-BIN | 2929 | Number of SVs per 1-Mbp bin, and per chromosome, normalised against the total number of SV per sample |
| | INDEL-BIN | 2757 | Number of SNVs per 1-Mbp bin, and per chromosome, normalised against the total number of INDEL per sample |
| Mutation type | MUT-WGS | 150 | Type of single-nucleotide substitution, double- and triple-nucleotide substitution (plus its adjacent nucleotide neighbours) |
| Driver gene/pathway | GEN | 554 | Presence of an impactful mutation in a suspected driver gene |
| | MOD | 1865 | Presence of an impactful mutation in a gene belonging to a suspected driver pathway |

# RF for individual features



**Fig. 1 Comparison of tumour-type classifiers using single and multiple feature types. a** Radar plots describing the cross-validation-derived accuracy (F1) score of Random Forest classifiers trained on each of 7 individual feature categories, across six representative tumour types. **b** Summary of Random Forest classifier accuracy (F1) trained on individual feature categories across all 24 tumour types. **c** Accuracy of classifiers trained on multiple feature categories. *RF Best Models* corresponds to the cross-validation F1 scores of Random Forest classifiers trained on the three best single-feature categories for all 24 tumour types. *DNN Model* shows the distribution of F1 scores for held-out samples for a multi-class neural network trained using passenger mutation distribution and type. *DNN Model + Drivers* shows F1 scores for the neural net when driver genes and pathways are added to the training features. The centre line in the boxplot represents the median of the F1 scores. The lower and upper bounds of the box represent the first and third quartile. The whiskers extend to 1.5 IQR plus the third quartile or minus the first quantile.
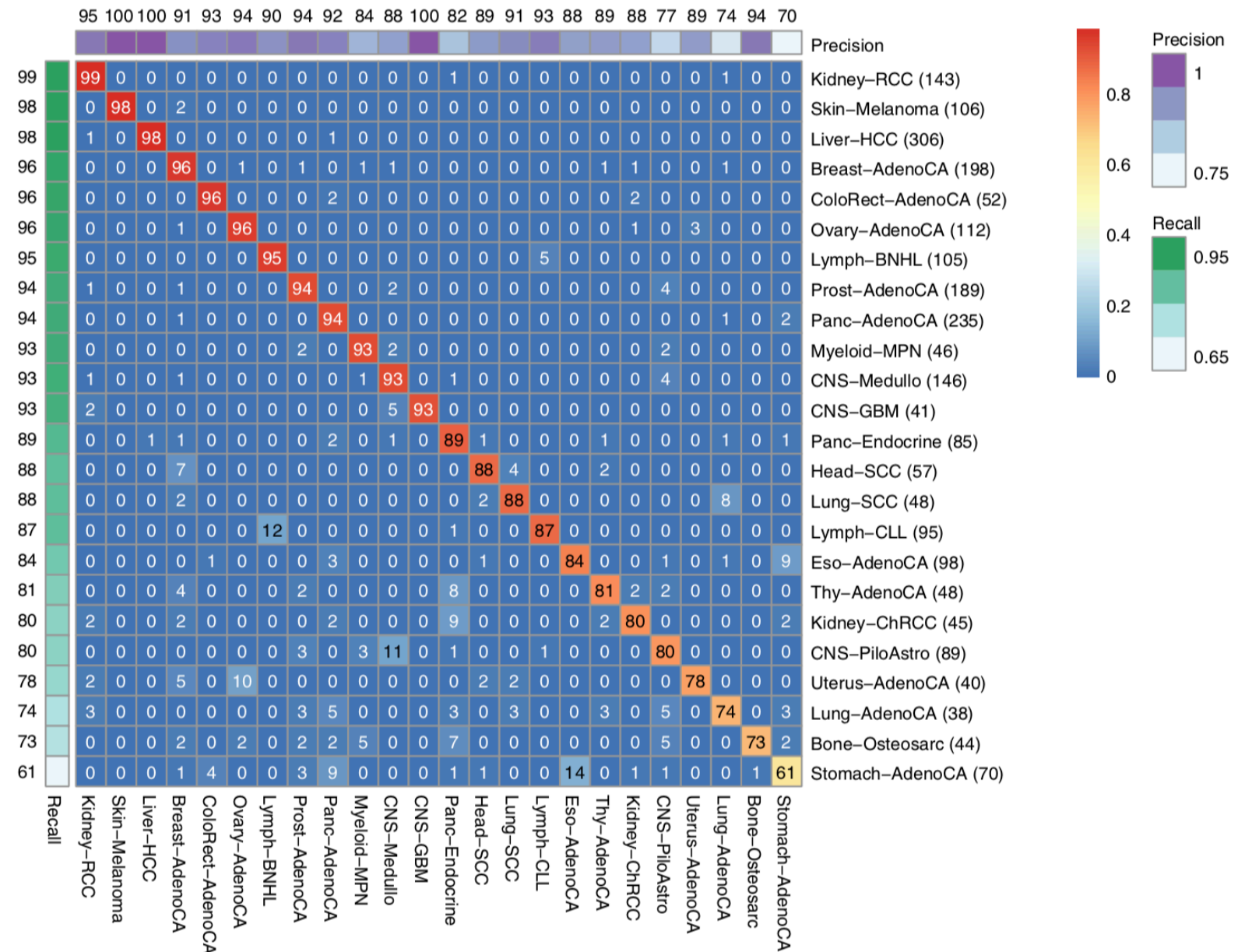
# Results (1/4)



**Fig. 2 Heatmap displaying the accuracy of the merged classifier using a held-out portion of the PCAWG data set for evaluation.** Each row corresponds to the true tumour type; columns correspond to the class predictions emitted by the DNN. Cells are labelled with the percentage of tumours of a particular type that were classified by the DNN as a particular type. The recall and precision of each classifier are shown in the colour bars at the top and left sides of the matrix. All values represent the mean of 10 runs using selected data set partitions. Due to rounding of values, some rows add up to slightly more or less than 100%.
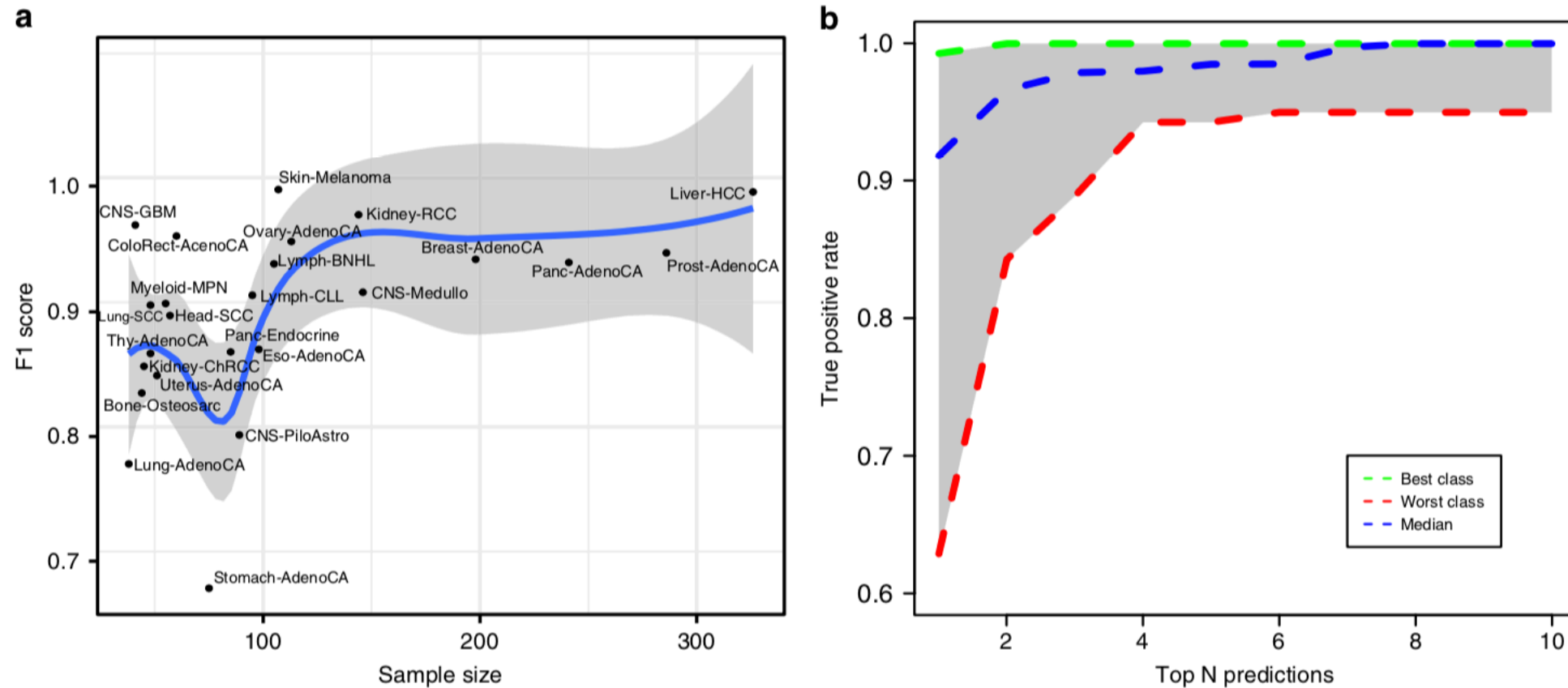
# Results (2/4)



**Fig. 3 Performance of the DNN on held-out PCAWG data. a** The relationship between training set size and prediction accuracy of the DNN is shown for each tumour type. The blue line represents a regression line fit using LOESS regression, while the grey area represents a 95% confidence interval for the regression function. **b** Accuracy of the classifier when it is asked to identify the correct tumour type among its top N-ranked predictions. The blue dashed line is the median true-positive rate among all 24 tumour classes. The green and red dashed lines correspond to the true-positive rate for the best- and worst-performing tumour classes.

# Results (3/4)

we applied the classifier trained on PCAWG sam-ples to an independent validation set of 1436 cancer whole genomes assembled from a series of published non-PCAWG projects. The validation set spans 14 distinct tumour types assembled from 21 publications or databases (Supplementary Data 4)
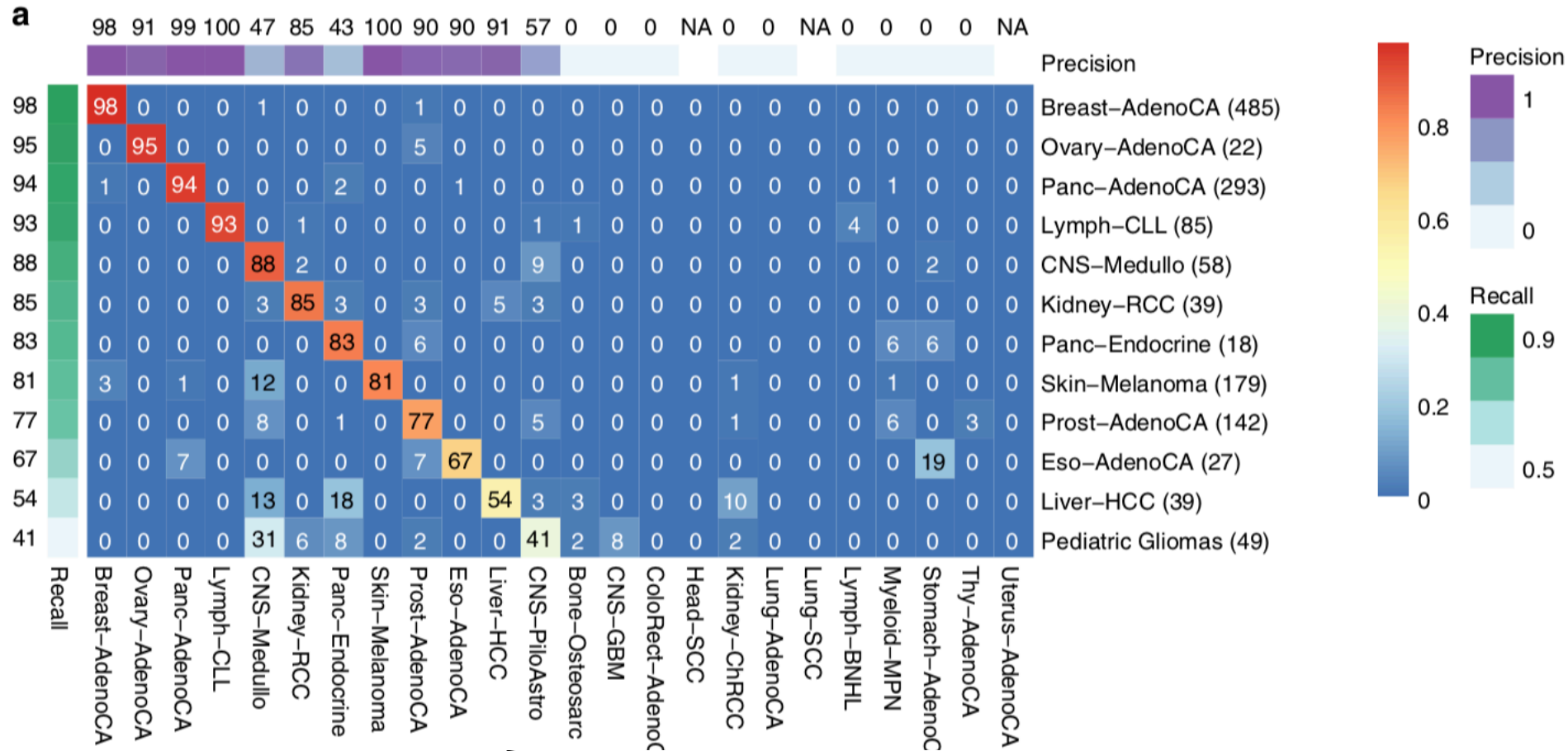


**Fig. 4 Prediction accuracy for the DNN against two independent validation data sets. a** Primary tumours. **b** Metastatic tumours. Each row corresponds to the true tumour type; columns correspond to the class predictions emitted by the DNN. Cells are labelled with the percentage of tumours of a particular type that were classified by the DNN as a particular type. The recall and precision of each classifier are shown in the colour bars at the top and left sides of the matrix. Due to rounding of values, some rows add up to slightly more or less than 100%.

# Results (4/4)

that combined a published series of 92 metastatic Panc-AdenoCA25 with an unpublished set of 2,028 metastatic tumours from known primaries across 16 tumour types recently sequenced by the Hartwig Medical Foundation (HMF), resulting in a combined set of 2120 samples across 16 tumour types (Sup- plementary Data 4).