

Defining the molecular state of human cancer

Biswajyoti Sahu^{1,2,3}, Päivi Pihlajamaa^{1,3}, Kaiyang Zhang⁴, Kimmo Palin^{1,2}, Saija Ahonen^{1,2},
Alejandra Cervera⁴, Ari Ristimäki^{1,2,5}, Lauri A. Aaltonen^{1,2}, Sampsa Hautaniemi⁴ and Jussi
Taipale^{1,3,6*}

1. *Applied Tumor Genomics Research Program, Faculty of Medicine, University of Helsinki, Helsinki, FI-00014, Finland*
2. *Medicum, Faculty of Medicine, University of Helsinki, Helsinki, FI-00014, Finland.*
3. *Department of Biochemistry, University of Cambridge, Cambridge, CB2 1GA, United Kingdom*
4. *Research Program in Systems Oncology, Faculty of Medicine, University of Helsinki, Helsinki, FI-00014, Finland*
5. *Department of Pathology, HUSLAB, Helsinki University Hospital, University of Helsinki, Helsinki, FI-00014, Finland*
6. *Department of Medical Biochemistry and Biophysics, Karolinska Institutet, Stockholm, SE-17177, Sweden*

* Correspondence: Jussi Taipale (ajt208@cam.ac.uk)

SUMMARY

Cancer is the most complex genetic disease known, with mutations implicated in more than 250 genes. However, it is still elusive which specific mutations found in human patients lead to tumorigenesis. Here we show that a combination of oncogenes that is characteristic of liver cancer (CTNNB1, TERT, MYC) induces senescence in human fibroblasts and primary hepatocytes. However, reprogramming fibroblasts to a liver progenitor fate, induced hepatocytes (iHeps), makes them sensitive to transformation by the same oncogenes. The transformed iHeps are highly proliferative, tumorigenic in nude mice, and bear gene expression signatures of liver cancer. These results show that tumorigenesis is triggered by a combination of three elements: the set of driver mutations, the cellular lineage, and the state of differentiation of the cells along the lineage. Our results indicate that cell identity is a key determinant in transformation, and establish a paradigm for defining the molecular states of human cancer.

Keywords: oncogenes, lineage-specific transcription factors, liver cancer, transdifferentiation, transformation

INTRODUCTION

Cancer genetics and genomics have identified a large number of genes implicated in human cancer (Forbes et al., 2017; Garraway and Lander, 2013; Vogelstein et al., 2013). Although some genes such as *p53* and *PTEN* are commonly mutated in many different types of cancer, most cancer genes are more lineage-specific. It is well established that human cells are harder to transform than rodent cells (Boehm et al., 2005; Chaffer and Weinberg, 2015; Kamijo et al., 1997; Metz et al., 1995; Rangarajan et al., 2004; Ruley, 1983; Stevenson and Volsky, 1986), which can be transformed using only MYC and RAS oncogenes (Land et al.,

1983; Shih et al., 1981; Sinn et al., 1987). Seminal experiments by Hahn and Weinberg established already 20 years ago that different human cell types can be transformed using a set of oncogenes that includes the powerful viral large-T and small-T oncoproteins from the SV40 virus (Hahn et al., 1999). Despite this early major advance, determining which specific mutations found in human patients lead to tumorigenesis has proven to be exceptionally difficult. This is because although viral oncoproteins are linked to several cancer types (Moore and Chang, 2010), most major forms of human cancer result from mutations affecting tumor-type specific sets of endogenous proto-oncogenes and tumor-suppressors (Haigis et al., 2019). The fact that the combinations of oncogenes are distinct between tumor types suggests that cell lineage-specific factors could somehow interact with oncogenes to drive most cases of human cancer, confounding mechanistic studies utilizing simple model cell types. This prompted us to systematically investigate the factors required for transformation of human cells using a combination of cell fate conversion and oncogene activation.

RESULTS

Generating proliferative induced hepatocytes using defined transcription factors and oncogenic drivers

Many human cell types can be converted to other cell types via a pluripotent state (Takahashi et al., 2007). However, as pluripotent cells are tumorigenic in nude mice, we chose to use direct lineage conversion (Davis et al., 1987; Pang et al., 2011; Sekiya and Suzuki, 2011) in combination with oncogene expression to identify the set of factors that define a particular type of human cancer cell. For this purpose, we developed a cellular transformation assay protocol, in which human fibroblasts (HF) are converted to induced hepatocytes (iHeps) using lentiviral overexpression of a combination of lineage-specific transcription factors (TF), followed by ectopic expression of liver cancer-specific oncogenes

(Figure 1A). Transdifferentiation of fibroblasts to iHeps has previously been reported by several groups (Du et al., 2014; Huang et al., 2014; Morris et al., 2014; Sekiya and Suzuki, 2011). To identify an optimal protocol for generating iHeps from HF_s (from human foreskin), we tested the previously reported combinations of TFs in parallel transdifferentiation experiments and analyzed the efficiency of iHep conversion by measuring the mRNA levels for liver markers (Du et al., 2014; Huang et al., 2014; Morris et al., 2014) such as *ALBUMIN*, *TRANSFERRIN*, and *SERPINA1* at different time points (**Figure 1B**, **Figure S1**). The combination of three TFs, HNF1A, HNF4A and FOXA3 (Huang et al., 2014) resulted in the most efficient iHep generation, based on the observation that out of all combinations tested, this combination resulted in the highest expression level of liver-specific genes at two, three, and four weeks after iHep induction (**Figure 1B**). This protocol also resulted in most efficient lineage conversion based on the analysis of cell morphology; by two weeks after iHep induction, the cells lost their fibroblast phenotype and formed spherical iHep progenitor colonies, from which immature, proliferative iHeps migrated outward. The iHeps fully matured to non-proliferative iHeps by six to seven weeks after induction (**Figure 2A**, **Figure S2**).

Oncogene exposure transforms induced hepatocytes but not control fibroblasts

To determine whether the iHeps could be transformed to liver cancer-like cells, we first plated immature (1 to 3 weeks post-transdifferentiation) iHeps on collagen-coated dishes and maintained them in hepatocyte culture media (HCM). Under such conditions, the iHeps mature, and their proliferation is arrested after two to three passages (Huang et al., 2014); after this point, further passaging induces cell death (**Figure 2B**). To confer the immature iHeps with unlimited proliferative potential and to drive them towards tumorigenesis, we transduced them with a set of the most common driver genes for liver cancer using lentiviral

constructs. For this purpose, we chose the five oncogenic drivers with the highest number of recurrent genetic alterations reported for liver cancer or hepatocellular carcinoma (HCC; from COSMIC, <https://cancer.sanger.ac.uk/cosmic>); these included four oncogenes, telomerase (*TERT*), β -catenin (*CTNNB1*), PI3 kinase (*PIK3CA*), and the transcription factor NRF2 (*NFE2L2*), as well as one tumor suppressor, p53 (*TP53*). In addition, we included the oncogene *MYC*, which is under tight control in normal cells (Lowe et al., 2004), but overexpressed in many cancer types, including HCC (Kalkat et al., 2017). Lentiviral expression of the fluorescent reporter mCherry with the oncogenic drivers in different combinations revealed that the pool of three oncogenes, *i.e.* constitutively active β -catenin (*CTNNB1*^{T41A}), *MYC* and *TERT*, together with *TP53* inactivation by CRISPR-Cas9 (*CMT*+*sgTP53*) resulted in highly proliferative iHeps with apparently unlimited proliferative potential (> 50 passages over more than one year; **Figure 2B**). Importantly, expression of the three oncogenes *CTNNB1*^{T41A}, *MYC* and *TERT* (*CMT*) alone also resulted in similar iHeps with long-term proliferative potential (**Figure 2B**). By contrast, ectopic expression of these oncogenic drivers in HF^s failed to yield transformed, proliferating fibroblasts, and rather resulted in cellular senescence and loss of the oncogene-transduced cells from the fibroblast population (**Figure 2B**). This is the first instance to our knowledge where HF^s can be directly transformed using this minimal combination of defined factors, indicating that lineage-specific TFs are the missing link for human cellular transformation using oncogenic drivers.

Tumorigenic properties of the transformed iHeps

To test for the tumorigenicity of the proliferative iHeps, we performed xenograft experiments. Subcutaneous injection of the *CMT*+*sgTP53* transformed iHeps, but not the control fibroblasts or iHeps with lineage-specific TFs alone, into nude mice resulted in tumor formation (**Figures 3A and 3B**). The process was reproducible in subsequent experiments; in

addition, the effect was not specific to the fibroblast line used, as we also successfully reprogrammed another HF cell line (human fetal lung fibroblast) using the same lineage-specific TFs, and transformed it using the same set of oncogenic drivers. The xenograft tumors from the CMT+sg*TP53* transformed iHeps derived from either fibroblast line can be detected by *in vivo* fluorescent imaging as early as 11-12 weeks (**Figure 3B**). Importantly, the histology of CMT+sg*TP53* tumors harvested at 20 weeks show highly malignant and proliferative HCC-like features (**Figure 3C**). Similarly, the CMT-transformed iHeps without *TP53* inactivation also resulted in tumor formation in nude mice 12 weeks post-injection (**Figure 3B**). These results demonstrate that both CMT and CMT+sg*TP53* transformed iHeps are tumorigenic, and indicate that ectopic expression of defined lineage-specific TFs and oncogenes can reprogram and transform HFs into cells that can robustly initiate tumors in nude mice.

Cancer genomes harbor large-scale chromosomal aberrations and are characterized by aneuploidy (Palin et al., 2018; Taylor et al., 2018). To understand the gross chromosomal aberrations in the transformed tumorigenic CMT and CMT+sg*TP53* iHeps compared to normal HFs, we performed spectral karyotyping, which showed a normal diploid male (46, XY) in HFs and aneuploid karyotypes in transformed iHeps (**Figure 3D**). The aneuploid transformed iHeps with CMT+sg*TP53* at early passage were characterized by two different populations with two distinct modal chromosome numbers (**Figure 3D**). The modal chromosome number of the first population was 45, XY, whereas the second population was pseudotetraploid, with a modal chromosome number between 67-92, XY; this pseudotetraploid state was consistently observed in late passage transformed iHeps. The major chromosomal aberrations that were similar between the two populations were missing copies of chromosomes 4 and 13, a derivative of chromosome 19 containing a small portion of chromosome 3 [t3:19], an extra copy of Y and a loss of most of the p arm of chromosome

2. In comparison, the most common chromosomal aberrations reported in HCC are the gains of 1q (suggested target genes include *WNT14*, *FASL*) and 8q (*MYC*, *WISP1*) and the loss of 17p (*TP53*, *HIC1*), followed by losses of 4q (*LEF1*, *CCNA*) and 13q (*RBI*, *BRCA3*) (Cancer Genome Atlas Research Network, 2017; Moinzadeh et al., 2005) (**Figure 3E**). The first three chromosomal aberrations are expected not to be present in our case, as the transformation protocol leads to activation of the Wnt pathway and MYC expression, and loss of p53. Consistently with this, we did not observe lesions in 1q, 8q or 17p in our cells. However, other common aberrations found in HCC cells, loss of chromosomes 4 and 13 were detected in our transformed CMT+sg*TP53* iHep cells (**Figures 3D and 3E**). However, these chromosomal aberrations appeared not to be necessary for formation of tumors, as in the absence of targeted loss of p53 in CMT iHep cells, we did not observe these lesions (**Figure 3D**). However, both CMT+sg*TP53* and CMT iHeps displayed pseudotetraploidy, similar to what is commonly observed in HCC (**Figures 3D and 3E**). These results indicate that the transformed iHeps have similar chromosomal aberrations to those reported earlier in liver cancer, consistent with their identity as HCC-like cells.

Transformed iHeps show gene expression profile similar to liver cancer

To understand the gene expression dynamics and to map the early events of lineage conversion and oncogenic transformation, we performed single cell RNA-sequencing (scRNA-seq) of HFs and iHeps with or without oncogene transduction. The cells were clustered according to their expression profiles using Seurat (Satija et al., 2015); a total of 15 separate clusters of cells were identified during the course of the transdifferentiation and reprogramming and visualized by t-distributed stochastic neighbour embedding (t-SNE) plots (van der Maaten and Hinton, 2008) (**Figures 4A and 4B**). Importantly, the scRNA-seq

indicated that the CMT-transformed iHeps are a clearly distinct population of cells compared to the iHeps, whereas CMT-transduced HF cells are more similar to the control HF cells (**Figure 4B**).

To determine the trajectory of differentiation of the cells, we performed RNA velocity analysis (La Manno et al., 2018), which determines the direction of differentiation of individual cells based on comparison of levels of spliced mRNAs (current state) with nascent unspliced mRNAs (representative of future state). This analysis confirmed that the cell populations analyzed were differentiating along the fibroblasts–iHep–transformed iHep axis (**Figure 4C**). We next identified marker genes for each cell cluster (see **Methods**). This analysis revealed that CMT-iHeps have a distinct gene expression signature and that they have lost the fibroblast gene expression program during the course of the reprogramming (**Figure S3A**). These results indicate that the iHep conversion and transformation have led to generation of liver-cell like transformed cells.

To further analyze gene expression changes during reprogramming and transformation, we performed pseudo-temporal ordering analysis of the scRNA-seq (see **Methods**). Consistently with the RNA velocity analysis, the pseudotime analysis showed transition from fibroblasts to iHeps and subsequently to CMT-transformed iHeps (**Figure S3B**). Similarly, CMT-transduced HF cells were ordered across pseudo-temporal time line (**Figure S3B**). The scRNA-seq analyses allow detection of the precise early events that occur during iHep formation and the origin of HCC by mapping the gene expression changes in the cells across the pseudotime. Furthermore, analyzing the molecular changes upon CMT-transduction provided mechanistic understanding of why oncogenes fail to transform HF cells without iHep conversion (**Figure 4D**); the pseudotime analysis of gene expression changes from CMT-transduced HF cells at one and three weeks compared to control HF cells was highly similar to the previously reported signature of cellular senescence (17 out of 18 genes; Marthandan et al., 2016) (**Figure 4D**). The senescence signature was much weaker both

during transdifferentiation of the iHeps and during their transformation (**Figure S3C**); instead, during iHep differentiation, the expression of non-canonical Wnt pathway components, including Wnt5a ligand and the Frizzled 5 receptor, were upregulated (**Figure 4D**). During transformation, the exogenous CTNNB1^{T41A} activated the canonical Wnt pathway, suppressing expression of the non-canonical ligand Wnt5a. We also observe activation of the NOTCH pathway early during tumorigenesis; expression of *NOTCH1*, *NOTCH3* and their ligand *JAG1* (**Figure 4D**, top) are strongly upregulated, together with the canonical NOTCH target gene *HES1* (Borggrefe and Oswald, 2009) and the liver specific target *NR4A2* (Zhu et al., 2017). These results are consistent with the proposed role of the NOTCH pathway in liver tumorigenesis (Villanueva et al., 2012; Zhu et al., 2017). Thus, we show that CMT transduced HF's undergo senescence, whereas the proliferative phenotype of CMT-iHep cells is associated with gene expression changes in Wnt and NOTCH signaling pathways.

To determine whether the gene expression signatures observed in transformed iHeps were similar to those observed in human liver tumors, we compared the scRNA-seq results to the published liver cancer data sets (Cancer Genome Atlas Research Network, 2017). Majority of the cases in TCGA HCC pan-cancer dataset show genetic alterations in CMT-iHep-specific marker genes (**Figure S3A**, 69% of 372 cancer cases) and the CMT-iHep-specific markers showed larger overlap with HCC when compared to cancers of pancreas and prostate, suggesting the specificity of this set of genes for liver tumorigenesis (**Figure S3D**). We also analyzed the expression of the CMT-iHep marker genes that show genetic alterations in TCGA liver cancer data across the pseudotime in our scRNA-seq data. The expression of this subset of the CMT-iHep marker genes was also clearly increased, lending further credence to the fact that upregulation of this set of genes is an early event in liver tumorigenesis (**Figure S3E**).

Direct conversion of human fibroblasts to liver cancer cells results in up-regulation of liver cancer markers

To determine the changes in gene expression and chromatin accessibility in the proliferative iHeps, we first performed bulk RNA-seq analysis from the tumorigenic CMT and CMT+sg*TP53* iHeps that were used for the xenograft implantation, as well as cells derived from the resulting tumors. Importantly, the genes that were differentially expressed in both CMT- and CMT+sg*TP53*-transformed iHeps compared to fibroblasts showed a clear and significant positive enrichment for the previously reported “subclass 2” liver cancer signature (Hoshida et al., 2009), associated with proliferation and activation of the MYC and AKT signaling pathways (**Figure 4E**). The effect was specific to liver cancer, as we did not observe significant enrichment of gene expression signatures of other cancer types (**Figure S4**). During the reprogramming, we observed a clear up-regulation of common liver marker genes such as *ALB*, *APOA2*, *SERPINA1*, and *TF*, and down-regulation of fibroblast markers such as *MMP3*, *FGF7*, *THY1*, and *FAP*, in proliferative and tumorigenic iHeps. Importantly, the xenograft tumor from the CMT+sg*TP53* cells retained similar liver-specific gene expression profile (**Figure 5A**). We also detected a clear up-regulation of several liver cancer marker genes such as *AFP*, *GPC3*, *SAA1*, and *VILI* in transformed iHeps and in CMT+sg*TP53* tumors compared to control fibroblasts (**Figure 5A**); *AFP* was also found among the most enriched genes (**Figure S5A**) in both CMT+sg*TP53*- and CMT-transformed iHeps. Furthermore, we observed a negative correlation between the CMT+sg*TP53* and CMT iHep specific genes and the genes positively associated with liver cancer survival (**Figure S5B**), lending further credence to liver cancer-identity of the CMT+sg*TP53* and CMT transformed iHeps.

ATAC-seq analysis of the fibroblasts and CMT+sg*TP53* cells revealed that the changes in marker gene expression were accompanied with robust changes in chromatin accessibility at the corresponding loci (**Figure 5B**). To assess chromatin accessibility and DNA methylation at a single-allele level, we performed NaNoMe-seq (see **Methods**), where accessible chromatin is methylated at GpC dinucleotides using the bacterial methylase M.CviPI (Kelly et al., 2012). Sequencing of the genome of the treated cells using single-molecule Nanopore sequencer then allows both detection of chromatin accessibility (based on the presence of methylated cytosines at GC dinucleotides) and DNA methylation at CG dinucleotides. This analysis confirmed the changes in DNA accessibility detected using ATAC-seq (**Figure 5C**). Changes in DNA methylation at promoters of the differentially expressed genes were relatively minor (**Figure 5C**), suggesting that the mechanism of reprogramming does not critically depend on changes in CpG methylation at the marker loci. Taken together, these results indicate that our novel cell transformation assay using lineage-specific TFs and cancer-specific oncogenes can reprogram fibroblasts to lineage-specific cancer that bears a gene expression signature similar to that observed in HCC.

Cellular lineage and the differentiated state of cells along the lineage are critical for tumorigenesis

To identify the necessary and sufficient factors that define lineage-specific cancer types we have here developed a novel cellular transformation protocol, and, for the first time, report direct conversion of HFs to liver cancer cells. First, lentiviral overexpression of three lineage-specific TFs reprograms HFs to iHeps, and subsequent ectopic expression of liver cancer-specific oncogenic factors transforms iHeps to a highly proliferative and tumorigenic phenotype with chromosomal aberrations and gene expression signature patterns similar to HCC. Importantly, lineage-conversion by specific TFs is required for the transformation

process since the same oncogenic drivers alone do not transform HF_s (**Figure 6A**). After lineage conversion by the defined TF_s, oncogenes alone (MYC, CTNNB1 and TERT) are sufficient to drive the transformation with or without inactivation of the tumor suppressor *TP53*. In contrast, oncogene transduction induces senescence in both HF_s and in differentiated adult human hepatocytes (**Figure 6A**). Interestingly, ectopic expression of TF_s used for iHep conversion (HNF1A, HNF4A, and FOXA3) one week prior to oncogene transduction protects the hepatocytes from oncogene-induced senescence (**Figure 6A**). These results show that fully differentiated hepatocytes are not susceptible for transformation by the liver-specific set of oncogenes, indicating that in addition to cellular lineage, also the differentiated state of cells along the lineage is critical for tumorigenesis. This finding is consistent with our experiments studying transformation in reprogrammed induced neurons (iNs). As these cells become terminally differentiated and post-mitotic immediately upon transdifferentiation, neither medulloblastoma nor neuroblastoma-specific oncogenes were able to make them re-enter the cell cycle (**Figure 6A and Figure S6**). These results establish a paradigm for testing the tumorigenicity of combinations of cancer genes, and their interactions with cellular lineage and differentiated state (**Figure 6B**). In addition, reprogramming normal cells to cancer cells allow “live” analysis of the early stages of the tumorigenic program, facilitating approaches towards early molecular detection and prevention of cancer.

DISCUSSION

In the past half-century, a very large number of genetic and genomic studies have been conducted using increasingly powerful technologies, resulting in identification of more than 250 genes that are recurrently mutated in cancer. However, in most cases, the evidence that the mutations in the genes actually cause cancer is correlative in nature, and requires

assumptions about background mutation frequency and rates of clonal selection in normal tissues (Martincorena et al., 2018). Furthermore, cancer genes are known to act in combination, and determining candidate sets of genes that are sufficient to cause cancer using genetic data alone would require astronomical sample sizes. Mechanistic studies are thus critical for conclusively determining that a particular gene is essential for cancer formation, and for identification of sets of genes that are sufficient for tumorigenesis. In principle, individual driver genes and their combinations could be identified and validated using particular primary cell types. However, approaches using primary cells are severely limited by the fact that for most tissues, sufficient amounts of live human tissue material are hard to obtain. Furthermore, the cell type of origin for most cancer types is not known, and it is commonly assumed that tumors originate from rare and hard-to-isolate subpopulations of cells (e.g. stem cells, or transient progenitor cells in the case of pediatric tumors). Although some previous studies have reported oncogene combinations that can transform primary human cells (Drost et al., 2015; Matano et al., 2015; Park et al., 2018), many primary cells may not be at the specific differentiated state that is required for transformation (**Figure 6B**). For example, the combination of oncogenes required to transform cells to the most common prostate or lung cancer forms is not known. However, both primary prostate and lung epithelial cells can be transformed to small cell neuroendocrine carcinoma, suggesting that either a minority population of neuroendocrine cells are susceptible for transformation, or that both types of primary cells assume a neuroendocrine fate after transformation (Park et al., 2018).

Our approach allows more precise control of cell identity and differentiation state, facilitating analysis of interactions between driver genes, cell lineage and cell state (**Figure 6B**). Our results using the novel cellular transformation assay show that HFs can be directly converted to lineage-specific cancer. Using this assay, we were able to determine the

minimum events necessary for making human liver cancer in culture. By using lineage-specific TFs to generate the cell type of interest for transformation studies, our molecular approach can be generalized for identifying minimal determinants of any cancer type, paving the way towards elucidating the exact molecular mechanisms by which specific combinations of mutations cause particular types of human cancer.

ACKNOWLEDGMENTS

We thank Drs. Otto Kauko, Teemu Kivioja, and Minna Taipale for critical review of the manuscript, and Tomi Leung, Anu M. Luoto, Kaisu Jussila and Inga-Lill Åberg for technical assistance. We also thank HiLIFE research infrastructures including Biomedicum Virus Core, Single-cell Analytics FIMM, Biomedicum Imaging Unit and Laboratory Animal Center. This work was supported by grants from Academy of Finland (Finnish Center of Excellence Program; 2012-2017, 250345 and 2018-2025, 312041, Post-doctoral fellowships; 274555, 288836 and Research Fellowships, 317807), Jane and Aatos Erkko Foundation, and the Finnish Cancer Foundation.

AUTHOR CONTRIBUTIONS

Corresponding Author and Lead Contact: JT. JT conceived and supervised the study. BS designed and performed all the experiments with help from PP. BS performed the initial data processing of single cell and bulk RNA-seq data. KZ performed single cell RNA-seq analysis. BS performed the NaNoMe-seq and KP analyzed the data with inputs from SA and LA. AC and KZ performed bulk RNA-seq analysis with inputs from BS and PP. AR reviewed and analysed the immunohistochemistry slides, SH provided the bioinformatics support and inputs into the project. All authors contributed to the writing of the manuscript.

DECLARATION OF INTERESTS

The authors declare no competing interests.

FIGURE LEGENDS

Figure 1. Generating proliferative induced hepatocytes using defined transcription factors and oncogenic drivers.

(A) Schematic outline of the cell transformation assay for making lineage-specific cancer by lentiviral expression of three lineage-specific TFs to convert HFs to induced hepatocytes (iHep) and defined oncogenic drivers to transform iHeps to proliferating and tumorigenic cells.

(B) Comparison of TF combinations (Du et al., 2014; Huang et al., 2014; Morris et al., 2014) for converting human fibroblasts to iHeps by detecting transcript levels for liver marker genes (*ALBUMIN*, *TRANSFERRIN* and *SERPINA1/α-1-antitrypsin*) by qRT-PCR at different time points after iHep conversion, normalized to GAPDH levels (mean ± standard error).

Figure 2. Oncogene exposure transforms induced hepatocytes but not control fibroblasts.

(A) Phase contrast microscope images showing the phenotype and morphology of the cells in the course of conversion of fibroblasts to iHeps at different times points after transduction with a cocktail of three TFs HNF1A, HNF4A and FOXA3 (Huang et al., 2014).

(B) Generation of highly proliferative iHep cells by transducing iHeps with two pools of liver cancer-specific oncogenic drivers, a list of xenograft experiments in nude mice that were used to test the tumorigenicity of different conditions, and mutation rates of the oncogenic drivers as reported in the COSMIC database for HCC and MYC amplification as reported in Ref. (Cancer Genome Atlas Research Network, 2017). CMT pool contains three oncogenes

CTNNB1^{T41A}, MYC, and TERT, and CMT+sg*TP53* pool contains the same oncogenes along with constructs for *TP53* inactivation by CRISPR-Cas9. Phase contrast microscope images showing the phenotype and morphology of the cells. Oncogenes are co-transduced with fluorescent reporter mCherry for detection of transduced cells. Oncogene transduction to fibroblasts fails to transform the cells, passaging of oncogene-expressing fibroblasts results in cellular senescence as demonstrated by beta-galactosidase staining and loss of mCherry-positive oncogene-expressing cells from the fibroblast population. Passaging of iHeps without oncogenes results in apoptosis after few passages. Scale bar 1000 μ m unless otherwise specified.

Figure 3. Tumorigenic properties of the transformed iHeps.

(A) Subcutaneous injection of transformed iHeps results in xenograft tumors in nude mice (tumor size of 1.5 cm ~ 23 weeks after xenotransplantation). Proliferative iHeps transduced with defined CMT oncogenes with *TP53* inactivation (CMT+sg*TP53*) or control iHeps without oncogenes were used in the injections.

(B) *In vivo* imaging of xenograft tumors ~12 weeks after implantation. Two biological replicate experiments are shown for CMT+sg*TP53* cells with iHep conversion and oncogene transduction with *TP53* inactivation performed in two separate human fibroblast cell lines (foreskin fibroblast [*left panel*] and fetal lung fibroblast [*middle*]) as well as proliferative CMT iHeps without *TP53* inactivation (*right*). Fluorescence signal emitted by mCherry co-transduced with the oncogenes is detected *in vivo* using the Lago system (scale bar = radiance units). Control mice are injected with either fibroblasts or iHeps.

(C) Histological analysis of CMT+sg*TP53* tumor tissue harvested at 20 weeks. Hematoxylin-eosin (H&E) staining for general histology and immunohistochemical staining for Ki-67 for cell proliferation (100x magnification).

(D) Analysis of chromosomal aberrations in the transformed iHeps by spectral karyotyping. CMT+sg*TP53* cells were analyzed at passage 18 (*early*) and passage 50 (*late*) and CMT cells at passage 18. Fibroblasts have normal diploid karyotype (46, XY, representative spectral image on *left*) and transformed iHeps show aneuploidies as indicated in the figure. Early passage CMT+sg*TP53* cells show two different populations with two distinct modal chromosome numbers (45, XY and 67-92, XY, representative spectral image for 45, XY on *middle-left*). Late passage CMT+sg*TP53* cells have modal chromosome number 67-92, XY (*middle-right*) and CMT cells 75, XY (*right*).

(E) Frequencies of chromosomal alterations reported for human HCC samples (see Moinzadeh et al., 2005).

Figure 4. Transformed iHeps show gene expression profile similar to liver cancer.

(A-B) t-SNE plots of 3,500 single cells from fibroblasts, iHeps at one–three weeks after iHep induction, iHeps transduced with CMT oncogenes at one week and harvested for scRNA-seq two weeks later, and fibroblasts transduced with CMT oncogenes and harvested at one and three weeks. Cells are colored by sample (**A**), and distinct clusters (**B**) based on their expression profiles with sample collection time points indicated.

(C) Principal component analysis (PCA) projection of single cells from control fibroblasts, iHeps at one–three weeks after iHep induction, and CMT-iHeps two weeks after oncogenes shown with velocity field with the observed states of the cells shown as circles and the extrapolated future states shown with arrows for the first two principal components. Cells are colored by cluster identities corresponding to Figure 3B.

(D) Relative expression of the genes from the Notch signaling pathway (*panel on the right*) across pseudotime in the single-cell RNA-seq data from control fibroblasts, iHeps at one–three weeks after iHep induction, and CMT-iHeps two weeks after oncogenes (the expression

of a gene in a particular cell relative to the average expression of that gene across all cells). Relative expression of the senescence marker genes (Marthandan et al., 2016) (*panel on the left*) from control fibroblasts and fibroblasts transduced with CMT oncogenes and harvested at one and three weeks after transduction. Color codes illustrating sample and cluster identities correspond to the colors in Figures 3A and 3B, respectively.

(E) Gene set enrichment analysis (GSEA) results for CMT-iHeps and CMT+sg*TP53*-iHeps compared to control fibroblasts against liver cancer signature (Subclass 2; Hoshida et al., 2009) from molecular signatures database (MSigDB). Positive normalized enrichment score (NES) reflects overrepresentation of liver cancer signature genes among the top ranked differentially expressed genes in CMT-iHep and CMT+sg*TP53*-iHep conditions compared to control fibroblasts.

Figure 5. Direct conversion of human fibroblasts to liver cancer cells results in up-regulation of liver cancer markers.

(A) Differential expression levels [$\log_2(\text{fold change})$] of marker genes for fibroblasts, hepatocytes, and liver cancer in bulk RNA-seq measurements from CMT+sg*TP53*-iHeps, CMT-iHeps and xenograft tumor from CMT+sg*TP53* cells against control fibroblasts (\pm standard error).

(B) IGV snapshots for promoter regions of representative genes from fibroblast markers (*MMP3*), liver markers (*SERPINA1*/ α -1-antitrypsin), and liver cancer markers (*SAAI*) showing ATAC-seq enrichment from fibroblast and CMT+sg*TP53*-iHeps.

(C) Chromatin accessibility and CpG methylation of DNA measured using NaNoMe-seq. Cytosine methylation detected using Nanopore sequencing from CMT+sg*TP53*-iHeps and control fibroblasts is shown for promoter regions of representative genes from fibroblast markers (*MMP3*), liver markers (*SERPINA1*/ α -1-antitrypsin), and liver cancer markers

(*SAAT*) using a window of TSS \pm 1500 bp. GpCpH methylation (all GC sequences where the C is not part of a CG sequence also, top) reports on chromatin accessibility, whereas HpCpG methylation reports on endogenous methylation of cytosines in the CpG context.

Figure 6. Cellular lineage and the differentiated state of cells along the lineage are critical for tumorigenesis.

(A) (*Top*) Beta-galactosidase staining as a marker of cellular senescence in primary human hepatocytes (control), after transduction of CMT oncogenes, or after transduction with iHep-TFs (HNF1A, HNF4A, FOXA3) followed by CMT oncogene transduction one week later (stained three weeks after first transduction). (*Middle*) Beta-galactosidase staining as a marker of cellular senescence in control fibroblasts and fibroblasts transduced with CMT oncogenes and stained at p4. (*Bottom*) Fluorescent microscope images of induced neurons with and without oncogene transduction (at three weeks of neuronal differentiation) visualized using EGFP at ten weeks after neuronal conversion.

(B) Schematic presentation of the molecular approach for identifying minimal determinants of tumorigenesis in specific tissues. Lineage-specific transcription factors are used to reprogram human fibroblasts to precise cellular identity (*left*), whose transformation by specific combinations of oncogenes (*right*) can then be tested. This approach combined with single-cell RNA-seq and RNA velocity analyses allows also analysis of which cell type along the stem cell to terminally differentiated cell axis (*top to bottom*) is susceptible for transformation.

METHODS

LEAD CONTACT AND MATERIALS AVAILABILITY

Further information and requests for resources and reagents should be directed to and will be fulfilled by the Lead Contact, Jussi Taipale (ajt208@cam.ac.uk).

EXPERIMENTAL MODELS AND SUBJECT DETAILS

Cell lines used in this study include human foreskin fibroblasts (HFF, CCD-1112Sk) and human fetal lung (HFL) fibroblasts that were obtained from ATCC (#CRL-2429 and #CCL-153, respectively) and cultured in fibroblast medium (DMEM supplemented with 10% FBS and antibiotics, Thermo Fisher Scientific). Primary adult human hepatocytes (#HUCPI, batch HUM4122A, Lonza) were plated on type I collagen-coated 24-well plates in plating medium (MP100, Lonza) and maintained in hepatocyte growth medium (HCM, Lonza) as per vendor's instructions. In xenograft experiments, 6-week old immunodeficient BALB/c nude male mice (Scanbur) were used.

METHODS DETAILS

Plasmids and lentiviral production

Full-length coding sequences for the TFs and oncogenes were obtained from GenScript and cloned into the lentiviral expression vector pLenti6/V5-DEST using the Gateway recombination system (Thermo Fisher Scientific). Expression construct for mCherry (#36084), lentiviral Cas9 expression construct LentiCas9-Blast (#52962), a cloning backbone lentiGuide-Puro (#52963), and the constructs for neuronal conversion (Tet-O-FUW-Ascl1, #27150; Tet-O-FUW-Brn2, #27151; Tet-O-FUW-Myt11, #27152; Tet-O-FUW-NeuroD1, # 30129; pTetO-Ngn2-Puro, #52047; Tet-O-FUW-EGFP, # 30130; FUW-M2rtTA, #20342) were obtained from Addgene. The six pairs of single-stranded oligos corresponding

to the guide sequences targeting the *TP53* gene in the GeCKO library were ordered from IDT, annealed, and ligated into lentiGuide-Puro backbone (Shalem et al., 2014). For virus production, the plasmids were co-transfected with the packaging plasmids psPAX2 and pMD2.G (Addgene #12260 and #12259, respectively) into 293FT cells (Thermo Fisher Scientific) with Lipofectamine 2000 (Thermo Fisher Scientific). Fresh culture medium was replenished on the following day, and the virus-containing medium was collected after 48 h. The lentiviral stocks were concentrated using Lenti-X concentrator (Clontech) and stored as single-use aliquots.

Generation of iHeps

LentiCas9-Blast virus was transduced to early-passage human fibroblasts (MOI = 1) with 8 µg/ml polybrene. Blasticidin selection (4 µg/ml) was started two days after transduction and continued for two weeks. Early passage blasticidin-resistant cells were used in the reprogramming experiments by transducing cells with constructs for TF expression in combinations reported earlier by Morris et al. (FOXA1, HNF4A, KLF5) (Morris et al., 2014), Du et al. (HNF4A, HNF1A, HNF6, ATF5, PROX1, CEBPA) (Du et al., 2014) and Huang et al. (FOXA3, HNF4A, HNF1A) (Huang et al., 2014) with MOI = 0.5 for each factor and 8 µg/ml polybrene (day 1). The medium was changed to fresh fibroblast medium containing β-mercaptoethanol on day 2 and to a defined hepatocyte growth medium (HCM, Lonza) on day 3. On day 6, the cells were passaged on plates coated with type I collagen (Sigma) in several technical replicates, and thereafter, the HCM was replenished every two–three days.

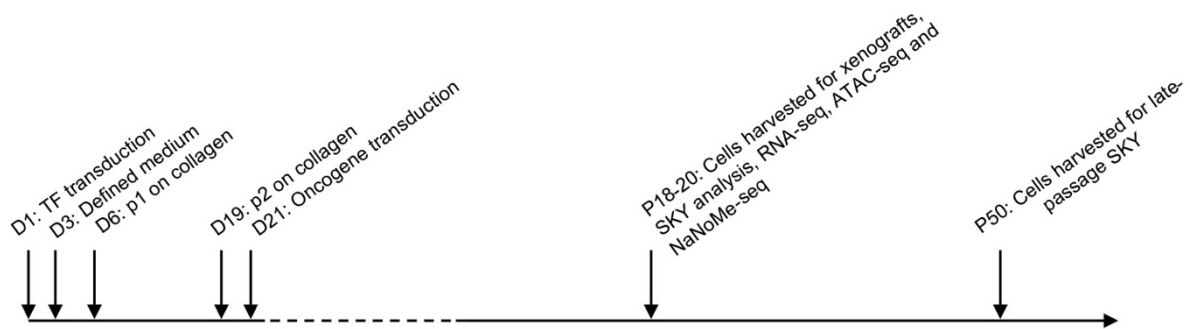
Primary adult human hepatocytes (#HUCPI, batch HUM4122A, Lonza) were plated on type I collagen-coated 24-well plates in plating medium (MP100, Lonza) and maintained in hepatocyte growth medium (HCM, Lonza) as per vendor's instructions. One day after plating, cells were transduced either with TFs that were used for iHep conversion (FOXA3,

HNF4A, HNF1A, MOI = 0.5) or with CMT oncogenes (CTNNB1^{T41A}, MYC, TERT, MOI = 1) with 8 µg/ml polybrene in HCM medium. Fresh HCM medium was replenished on the following day and regularly thereafter. Seven days after iHep-TF transduction, these cells were transduced with CMT oncogenes as above. Beta-galactosidase staining for senescence analysis was performed three weeks after plating the cells.

Generation of HCC-like cells

The iHeps generated using the three TFs (FOXA3, HNF4A, HNF1A) were passaged on type I collagen-coated plates on day 19 after iHep induction (p2) in HCM and transduced with different combinations of lentiviral constructs encoding the oncogenes (CTNNB1^{T41A}, MYC, TERT, NFE2L2, PIK3CA^{E545K}) on day 21 (MOI = 1 for each factor with 8 µg/ml polybrene). For CMT+sg*TP53* condition, the CMT oncogenes (CTNNB1^{T41A}, MYC and TERT) were transduced along with a pool of six sgRNAs targeting the *TP53* gene. Fresh HCM was replenished on the day following the transduction, cells were maintained in HCM, and passaged when close to confluent. From fifth passaging onwards after oncogene induction, cells were maintained in HCM supplemented with 1% defined FBS (Thermo Fisher Scientific). For single-cell RNA-sequencing experiments, the iHeps were transduced with CMT oncogenes (MOI = 1 with 8 µg/ml polybrene) on day 8 with fresh HCM replenished on day 9, and the cells were harvested for single-cell RNA-sequencing at the indicated time points from replicate culture wells. In all experiments, viral construct for mCherry expression was co-transduced with the oncogenes. As controls, HF^s were transduced with the same combination of oncogenes (CTNNB1^{T41A}, MYC, TERT, MOI = 1 for each factor with 8 µg/ml polybrene). Fresh medium was changed on the day following the transduction, and the cells were passaged regularly. Cells were harvested for scRNA-seq

analysis one and three weeks after transduction, and used for beta-galactosidase staining three weeks after transduction.



Generation of induced neurons (iN)

HF_s were plated on Matrigel-coated wells and transduced on the following day with tetracycline-inducible TF constructs for iN conversion (Ascl1, Brn2, Myt11, NeuroD1, and Ngn2) with MOI = 0.3 for each factor and 8 µg/ml polybrene with co-transduction of lentiviral construct for EGFP; day 1). The medium was changed to fresh fibroblast medium containing β-mercaptoethanol on day 2. Doxycycline induction (2 µg/ml) was started on day 5, and the medium was replaced with defined N2B27 neuronal medium supplemented with small molecules (CHIR, SB431542, LDN, dcAMP, and Noggin) and doxycycline on day 6 and cells were maintained in the defined medium thereafter. At three weeks of iN conversion, cells were transduced with one of the two oncogene pools specific either for neuroblastoma (ALK^{R1275Q}, MYCN, NRAS^{Q61R}, PIK3CA^{E545K}, BRAF^{V600E}, PTPN11, PDGFRA, KIT, IDH1^{R132H}) or for medulloblastoma (CTNNB1^{T41A}, NRAS^{Q61R}, PIK3CA^{E545K}, SMO^{W535L}, H3F3A^{K28M}, IDH1^{R132H}) with MOI = 0.5 for each factor and 8 µg/ml polybrene in neuronal medium. Fresh neuronal medium was replaced on the following day, the cells were maintained in neuronal medium and followed for 10-20 weeks.

Xenografts

Oncogene-induced CMT and CMT+sg*TP53* cells at p20 and control iHeps were harvested, 10^7 cells were resuspended in HCM supplemented with 1% defined FBS and mixed with equal volume of Matrigel (growth factor reduced basement membrane matrix, Corning #356231) and injected subcutaneously into the flank of a 6-week old immunodeficient BALB/c nude male mice (Scanbur). Similarly, 10^7 control fibroblasts were injected in equal volume of fibroblast medium and Matrigel. *In vivo* imaging of the tumors was performed for the mice under isoflurane anesthesia using the Lago system (Spectral Instruments Imaging). Photon counts from the mCherry were detected with fluorescence filters 570/630 nm and superimposed on a photographic image of the mice. Tumors were harvested 23-25 weeks after injection. All the experiments were performed according to the guidelines for animal experiments at the University of Helsinki and under license from appropriate Finnish Review Board for Animal Experiments.

SKY analysis

Spectral karyotype analysis was performed at Roswell Park Cancer Institute Pathology Resource Network. Cells were treated for 3 hours with 0.06 $\mu\text{g}/\text{ml}$ of colcemid, harvested and fixed with 3:1 methanol and acetic acid. Metaphase spreads from fixed cells were hybridized with SKY probe (Applied Spectral Imaging) for 36 hours at 37 degrees Celsius. Slides were prepared for imaging using CAD antibody kit (Applied Spectral Imaging) and counterstained with DAPI. Twenty metaphase spreads for each cell line were captured and analyzed using HiSKY software (Applied Spectral Imaging).

RNA isolation, qPCR and bulk RNA-sequencing

Total RNA was isolated from the control fibroblasts, iHeps harvested at day 5 and at weeks two, three, and four, CMT and CMT+sg*TP53* cells harvested at p20, and from tumor

tissues stored in RNALater (Qiagen), using RNeasy Mini kit (Qiagen) with on-column DNase I treatment. For qRT-PCR analysis, cDNA synthesis from two biological replicates was performed using the Transcriptor High-fidelity cDNA synthesis kit (Roche) and real-time PCR using SYBR green (Roche) with primers specific for each transcript (*GAPDH*-FP 5'-GGCCTCCAAGGAGTAAGACC, *GAPDH*-RP 5'-AGGGGAGATTCAGTGTGGTG, *ALB*-FP 5'-GGATGAAGGGAAGGCTTCGT, *ALB*-RP 5'-GAAATCTCTGGCTCAGGCGA, *TF*-FP 5'-GGCCACTAAGTGCCAGAGTT, *TF*-RP 5'-ATCCAGTGTCACAGCATCCG, *SERPINA1*-FP 5'-CTGTCTCCTCAGCTTCAGGC, *SERPINA1*-RP 5'-CACGAGACAGAAGACGGCAT). The Ct values for the target genes were normalized to those of GAPDH, and the mean values of sample replicates were shown for different conditions at the indicated time points. RNA-sequencing was performed from three biological replicate samples for each condition, using 400 ng of total RNA from each sample for poly(A) mRNA capture followed by stranded mRNA-seq library construction using KAPA stranded mRNA-seq kit for Illumina (Roche) as per manufacturer's instruction. Final libraries with different sample indices were pooled in equimolar ratios based on quantification using KAPA library quantification kit for Illumina platforms (Roche) and size analysis on Fragment Analyzer (AATI) and sequenced on HiSeq 4000 (Illumina).

For preprocessing and analysis of the RNA-Seq reads the SePIA pipeline (Icay et al., 2016) based on the Anduril framework (Ovaska et al., 2010) was used. Quality metrics from the raw reads were estimated with FastQC (<http://www.bioinformatics.babraham.ac.uk/projects/fastqc>) and Trimmomatic (Bolger et al., 2014) clipped adaptors and low-quality bases. After trimming, reads shorter than 20 bp were discarded. Kallisto (v0.44.0) with Ensembl v85 (Zerbino et al., 2018) was used for quantification followed by tximport (Soneson et al., 2015) and DESeq2 (Love et al., 2014) (v1.18.1) for differential expression calculating log₂(fold change) and standard error from

triplicate samples. Gene set enrichment analysis (Subramanian et al., 2005) was performed using GSEAPY (version 0.9.8) by ranking differentially expressed genes based on their $-\log_{10}(\text{p-value}) * \text{sign}(\text{fold-change})$ as metric. The gene signatures analysed for enrichment were collected from Molecular Signatures Database (MSigDB, version 6.2).

Single-cell RNA-sequencing

For single cell RNA-sequencing (scRNA-seq), iHeps and HF's at different time points were harvested, washed with PBS containing 0.04% bovine serum albumin (BSA), resuspended in PBS containing 0.04% BSA at the cell density of 1000 cells / μl and passed through 35 μm cell strainer. Library preparation for Single Cell 3'RNA-seq run on Chromium platform (10x Genomics) for 4000 cells was performed according to manufacturer's instructions and the libraries were paired-end sequenced (R1:27, i7-index:8, R2:98) on HiSeq 4000 (Illumina). Preprocessing of scRNA-seq data, including demultiplexing, alignment, filtering, barcode counting, and unique molecular identifier (UMI) counting was performed using CellRanger.

Quality control was applied separately for iHep and HFL-CMT samples. iHeps with fewer than 50,000 mapped reads, or expressing fewer than 4000 genes or with greater than 6% UMI originating from mitochondrial genes were excluded, while for HFL-CMT samples, cells with fewer than 2500 genes or with greater than 10% UMI originating from mitochondrial genes were excluded. All genes that were not detected in at least 5 cells were discarded. From each sample, 500 cells were down-sampled for further analysis. The data was normalized and log-transformed using Seurat (Satija et al., 2015) (version 3.0.2). A cell cycle phase-specific score was generated for each cell, across five phases (G1/S, S, G2/M, M and M/G1) based on Macosko et al. (Macosko et al., 2015) using averaged normalized expression levels of the markers for each phase. The cell cycle phase scores together with

nUMI and percentage of UMIs mapping to mitochondrial genes per cell were regressed out using a negative binomial model. The graph-based method from Seurat was used to cluster the cells. The first 30 PCs were used in construction of SNN graph, and 15 clusters were detected with a resolution of 0.8. Markers specific to each cluster were identified using the “negbinom” model. Pseudotime trajectories were constructed with URD (Farrell et al., 2018) (version 1.0.2). The RNA velocity analysis was performed using velocyto (La Manno et al., 2018) (version 0.17).

Oil-Red-O- and PAS-staining and beta-galactosidase activity assay

Oil-Red-O and Periodic Acid-Schiff (PAS) staining were performed according to the manufacturer’s recommendation (Sigma). Briefly, for Oil-Red-O-staining, cells were fixed with paraformaldehyde (4%) for 30 mins, washed with PBS, incubated with 60% isopropanol for 5 mins and Oil-Red-O working solution for 10 mins, and washed twice with 70% ethanol. For PAS-staining, cells were fixed with alcoholic formalin (3.7%) for 1 min, incubated with PAS solution for 5 mins and Schiff’s reagent for 15 mins with several washes with water between each step, and counter-stained with hematoxylin. Beta-galactosidase assay was performed using Senescence detection kit (Abcam) according to the manufacturer’s protocol for fixation and staining (overnight).

Immunohistochemistry

Tumor tissues were collected from mice injected with CMT+sgTP53 cells and fixed in 4% paraformaldehyde at 4°C overnight, dehydrated and embedded in paraffin. Five- μ m sections were stained with hematoxylin and eosin with Tissue-Tek DRS automated system at Tissue Preparation and Histochemistry Unit (University of Helsinki, Finland) using standard protocols. For Ki-67 detection, sections were dewaxed with xylene, rehydrated, boiled in 10

mM citrate buffer, and treated with 3% hydrogen peroxide for 5 min to block the endogenous peroxidase activity. Sections were incubated with Anti-Ki-67 antibody (sc-101861, SantaCruz) at 4°C overnight followed by 40 min incubation with Brightvision Poly-HRP-Anti Mouse staining reagent (ImmunoLogic). The immune complexes were visualized using DAB Quanto chromogen and substrate (ThermoFisher) and counterstained with hematoxylin. The slides were dehydrated and mounted using Eukitt (Sigma).

ATAC-seq

Fibroblasts and CMT+sg*TP53* cells (p20) were harvested and 50,000 cells for each condition were processed for ATAC-seq libraries using previously reported protocol (Corces et al., 2017) and sequenced PE 2x75 NextSeq 500 (Illumina). The quality metrics of the FASTQ files were checked using FASTQC and the adapters were removed using trim_galore. The reads were aligned to human genome (hg19) using BWA, and the duplicate reads and the mitochondrial reads were removed using PICARD. The filtered and aligned read files were used for peak calling using MACS2 and for visualizing the traces using the IGV genome browser.

NaNoME-seq (NOME-seq using Nanopore sequencing)

To profile chromatin accessibility using GC methylase using NOME-seq protocol (Kelly et al., 2012) and to utilize the ability of Nanopore sequencing to detect CpG methylation without bisulfite conversion and PCR, we adapted the NOME-seq protocol for Nanopore sequencing on Promethion (NaNoME-seq). The nuclei isolation and treatment with GC methylase (M.CviPI) was performed as described earlier (Kelly et al., 2012). The DNA was isolated from GC methylase treated nuclei by phenol chloroform followed by ethanol precipitation. The sequencing library for Promethion was prepared using the 1D genomic

DNA by ligation kit (SQK-LSK109) as per manufacturer's recommendation and we loaded 50 fmol of final adapter-ligated high molecular weight genomic DNA to the flow cells for sequencing. After sequencing and base calling, the Nanopore reads were aligned to GRCh37 reference genome with minimap2 (Li, 2018). Nanopolish (Simpson et al., 2017) was modified to call methylation in GC context. In total, 11Gbp of aligned read data from PCR amplified and GC methylated sequencing run was used to learn emission model for methylated GC sites. The learning process followed <https://github.com/jts/methylation-analysis/blob/master/pipeline.make> with adjustments for using human genome data and minimap2. For nuclear extract NaNoMe samples, methylation status was separately called for GC and CG sites. Similar independent method was recently described in a preprint by Lee et al (<https://www.biorxiv.org/content/10.1101/504993v2>). Reads with consecutive stretch of at least 80 GC sites with at least 75% methylated were filtered out due to expected cell free DNA contamination during library preparation as in Shipony et al. (<https://www.biorxiv.org/content/10.1101/504662v1>). The per site methylation levels in Figure 5C are mean smoothed with triangular kernel 5 sites wide. Fibroblast and CMT+sg*TP53* NaNoMe analyses used 20.3Gbp and 24.8Gbp of aligned data, respectively.

QUANTIFICATION AND STATISTICAL ANALYSIS

RNA-Seq: For preprocessing and analysis of the reads: SePIA pipeline based on the Anduril framework, clipping adaptors and low-quality bases: Trimmomatic. Kallisto (v0.44.0) with Ensembl v85 for quantification followed by tximport and DESeq2 (v1.18.1) for differential expression. Gene set enrichment analysis using GSEAPY (version 0.9.8).

Single cell RNA-seq: Preprocessing, demultiplexing, alignment, filtering, barcode counting, and unique molecular identifier (UMI) counting: CellRanger. Normalization and log-transformation: Seurat (version 2.3.4), cell cycle phase-specific scoring based on Macosko et

al. using averaged normalized expression levels of the markers for each phase. The cell cycle phase scores together with nUMI and percentage of UMIs mapping to mitochondrial genes per cell were regressed out using a negative binomial model. The graph-based clustering using Seurat. Markers specific to each cluster were identified using the “negbinom” model. Pseudotime trajectories were constructed with URD (version 1.0.2). The RNA velocity analysis was performed using velocity (version 0.17).

ATAC-seq: the adapters were removed using trim_galore, the reads were aligned to human genome (hg19) using BWA, and the duplicate reads and the mitochondrial reads were removed using PICARD. Peak calling using MACS2.

NaNoMe-seq: Nanopore reads were aligned using minimap2. Nanopolish was used to call GC methylation. For nuclear extract NaNoMe samples, methylation status was separately called for GC and CG sites. Similar independent method was recently described in a preprint by Lee et al (<https://www.biorxiv.org/content/10.1101/504993v2>). Reads with consecutive stretch of at least 80 GC sites with at least 75% methylated were filtered out due to expected cell free DNA contamination during library preparation as in Shipony et al.

(<https://www.biorxiv.org/content/10.1101/504662v1>).

DATA AND CODE AVAILABILITY

All sequence data is available under ENA accession PRJEB31262.

REFERENCES

Boehm, J.S., Hession, M.T., Bulmer, S.E., and Hahn, W.C. (2005). Transformation of human and murine fibroblasts without viral oncoproteins. *Mol Cell Biol* 25, 6464-6474.

Bolger, A.M., Lohse, M., and Usadel, B. (2014). Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* 30, 2114-2120.

Borggreffe, T., and Oswald, F. (2009). The Notch signaling pathway: transcriptional regulation at Notch target genes. *Cell Mol Life Sci* 66, 1631-1646.

Cancer Genome Atlas Research Network. (2017). Comprehensive and Integrative Genomic Characterization of Hepatocellular Carcinoma. *Cell* 169, 1327-1341 e1323.

Chaffer, C.L., and Weinberg, R.A. (2015). How does multistep tumorigenesis really proceed? *Cancer Discov* 5, 22-24.

Corces, M.R., Trevino, A.E., Hamilton, E.G., Greenside, P.G., Sinnott-Armstrong, N.A., Vesuna, S., Satpathy, A.T., Rubin, A.J., Montine, K.S., Wu, B., *et al.* (2017). An improved ATAC-seq protocol reduces background and enables interrogation of frozen tissues. *Nat Methods* 14, 959-962.

Davis, R.L., Weintraub, H., and Lassar, A.B. (1987). Expression of a single transfected cDNA converts fibroblasts to myoblasts. *Cell* 51, 987-1000.

Drost, J., van Jaarsveld, R.H., Ponsioen, B., Zimmerlin, C., van Boxtel, R., Buijs, A., Sachs, N., Overmeer, R.M., Offerhaus, G.J., Begthel, H., *et al.* (2015). Sequential cancer mutations in cultured human intestinal stem cells. *Nature* 521, 43-47.

Du, Y., Wang, J., Jia, J., Song, N., Xiang, C., Xu, J., Hou, Z., Su, X., Liu, B., Jiang, T., *et al.* (2014). Human hepatocytes with drug metabolic function induced from fibroblasts by lineage reprogramming. *Cell Stem Cell* *14*, 394-403.

Farrell, J.A., Wang, Y., Riesenfeld, S.J., Shekhar, K., Regev, A., and Schier, A.F. (2018). Single-cell reconstruction of developmental trajectories during zebrafish embryogenesis. *Science* *360*.

Forbes, S.A., Beare, D., Boutselakis, H., Bamford, S., Bindal, N., Tate, J., Cole, C.G., Ward, S., Dawson, E., Ponting, L., *et al.* (2017). COSMIC: somatic cancer genetics at high-resolution. *Nucleic Acids Res* *45*, D777-D783.

Garraway, L.A., and Lander, E.S. (2013). Lessons from the cancer genome. *Cell* *153*, 17-37.

Hahn, W.C., Counter, C.M., Lundberg, A.S., Beijersbergen, R.L., Brooks, M.W., and Weinberg, R.A. (1999). Creation of human tumour cells with defined genetic elements. *Nature* *400*, 464-468.

Haigis, K.M., Cichowski, K., and Elledge, S.J. (2019). Tissue-specificity in cancer: The rule, not the exception. *Science* *363*, 1150-1151.

Hoshida, Y., Nijman, S.M., Kobayashi, M., Chan, J.A., Brunet, J.P., Chiang, D.Y., Villanueva, A., Newell, P., Ikeda, K., Hashimoto, M., *et al.* (2009). Integrative transcriptome analysis reveals common molecular subclasses of human hepatocellular carcinoma. *Cancer Res* *69*, 7385-7392.

Huang, P., Zhang, L., Gao, Y., He, Z., Yao, D., Wu, Z., Cen, J., Chen, X., Liu, C., Hu, Y., *et al.* (2014). Direct reprogramming of human fibroblasts to functional and expandable hepatocytes. *Cell Stem Cell* *14*, 370-384.

- Icay, K., Chen, P., Cervera, A., Rantanen, V., Lehtonen, R., and Hautaniemi, S. (2016). SePIA: RNA and small RNA sequence processing, integration, and analysis. *BioData Min* *9*, 20.
- Kalkat, M., De Melo, J., Hickman, K.A., Lourenco, C., Redel, C., Resetca, D., Tamachi, A., Tu, W.B., and Penn, L.Z. (2017). MYC Deregulation in Primary Human Cancers. *Genes* (Basel) *8*.
- Kamijo, T., Zindy, F., Roussel, M.F., Quelle, D.E., Downing, J.R., Ashmun, R.A., Grosveld, G., and Sherr, C.J. (1997). Tumor suppression at the mouse INK4a locus mediated by the alternative reading frame product p19ARF. *Cell* *91*, 649-659.
- Kelly, T.K., Liu, Y., Lay, F.D., Liang, G., Berman, B.P., and Jones, P.A. (2012). Genome-wide mapping of nucleosome positioning and DNA methylation within individual DNA molecules. *Genome Res* *22*, 2497-2506.
- La Manno, G., Soldatov, R., Zeisel, A., Braun, E., Hochgerner, H., Petukhov, V., Lidschreiber, K., Kastrioti, M.E., Lonnerberg, P., Furlan, A., *et al.* (2018). RNA velocity of single cells. *Nature* *560*, 494-498.
- Land, H., Parada, L.F., and Weinberg, R.A. (1983). Tumorigenic conversion of primary embryo fibroblasts requires at least two cooperating oncogenes. *Nature* *304*, 596-602.
- Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* *34*, 3094-3100.
- Love, M.I., Huber, W., and Anders, S. (2014). Moderated estimation of fold change and dispersion for RNA-seq data with DESeq2. *Genome Biol* *15*, 550.

Lowe, S.W., Cepero, E., and Evan, G. (2004). Intrinsic tumour suppression. *Nature* *432*, 307-315.

Macosko, E.Z., Basu, A., Satija, R., Nemesh, J., Shekhar, K., Goldman, M., Tirosh, I., Bialas, A.R., Kamitaki, N., Martersteck, E.M., *et al.* (2015). Highly Parallel Genome-wide Expression Profiling of Individual Cells Using Nanoliter Droplets. *Cell* *161*, 1202-1214.

Marthandan, S., Baumgart, M., Priebe, S., Groth, M., Schaer, J., Kaether, C., Guthke, R., Cellerino, A., Platzer, M., Diekmann, S., *et al.* (2016). Conserved Senescence Associated Genes and Pathways in Primary Human Fibroblasts Detected by RNA-Seq. *PLoS One* *11*, e0154531.

Martincorena, I., Fowler, J.C., Wabik, A., Lawson, A.R.J., Abascal, F., Hall, M.W.J., Cagan, A., Murai, K., Mahbubani, K., Stratton, M.R., *et al.* (2018). Somatic mutant clones colonize the human esophagus with age. *Science* *362*, 911-917.

Matano, M., Date, S., Shimokawa, M., Takano, A., Fujii, M., Ohta, Y., Watanabe, T., Kanai, T., and Sato, T. (2015). Modeling colorectal cancer using CRISPR-Cas9-mediated engineering of human intestinal organoids. *Nat Med* *21*, 256-262.

Metz, T., Harris, A.W., and Adams, J.M. (1995). Absence of p53 allows direct immortalization of hematopoietic cells by the myc and raf oncogenes. *Cell* *82*, 29-36.

Moinzadeh, P., Breuhahn, K., Stutzer, H., and Schirmacher, P. (2005). Chromosome alterations in human hepatocellular carcinomas correlate with aetiology and histological grade--results of an explorative CGH meta-analysis. *Br J Cancer* *92*, 935-941.

Moore, P.S., and Chang, Y. (2010). Why do viruses cause cancer? Highlights of the first century of human tumour virology. *Nat Rev Cancer* *10*, 878-889.

Morris, S.A., Cahan, P., Li, H., Zhao, A.M., San Roman, A.K., Shivdasani, R.A., Collins, J.J., and Daley, G.Q. (2014). Dissecting engineered cell types and enhancing cell fate conversion via CellNet. *Cell* 158, 889-902.

Ovaska, K., Laakso, M., Haapa-Paananen, S., Louhimo, R., Chen, P., Aittomaki, V., Valo, E., Nunez-Fontarnau, J., Rantanen, V., Karinen, S., *et al.* (2010). Large-scale data integration framework provides a comprehensive view on glioblastoma multiforme. *Genome Med* 2, 65.

Palin, K., Pitkanen, E., Turunen, M., Sahu, B., Pihlajamaa, P., Kivioja, T., Kaasinen, E., Valimaki, N., Hanninen, U.A., Cajuso, T., *et al.* (2018). Contribution of allelic imbalance to colorectal cancer. *Nat Commun* 9, 3664.

Pang, Z.P., Yang, N., Vierbuchen, T., Ostermeier, A., Fuentes, D.R., Yang, T.Q., Citri, A., Sebastiano, V., Marro, S., Sudhof, T.C., *et al.* (2011). Induction of human neuronal cells by defined transcription factors. *Nature* 476, 220-223.

Park, J.W., Lee, J.K., Sheu, K.M., Wang, L., Balanis, N.G., Nguyen, K., Smith, B.A., Cheng, C., Tsai, B.L., Cheng, D., *et al.* (2018). Reprogramming normal human epithelial tissues to a common, lethal neuroendocrine cancer lineage. *Science* 362, 91-95.

Rangarajan, A., Hong, S.J., Gifford, A., and Weinberg, R.A. (2004). Species- and cell type-specific requirements for cellular transformation. *Cancer Cell* 6, 171-183.

Ruley, H.E. (1983). Adenovirus early region 1A enables viral and cellular transforming genes to transform primary cells in culture. *Nature* 304, 602-606.

Satija, R., Farrell, J.A., Gennert, D., Schier, A.F., and Regev, A. (2015). Spatial reconstruction of single-cell gene expression data. *Nat Biotechnol* 33, 495-502.

Sekiya, S., and Suzuki, A. (2011). Direct conversion of mouse fibroblasts to hepatocyte-like cells by defined factors. *Nature* 475, 390-393.

Shalem, O., Sanjana, N.E., Hartenian, E., Shi, X., Scott, D.A., Mikkelsen, T., Heckl, D., Ebert, B.L., Root, D.E., Doench, J.G., *et al.* (2014). Genome-scale CRISPR-Cas9 knockout screening in human cells. *Science* 343, 84-87.

Shih, C., Padhy, L.C., Murray, M., and Weinberg, R.A. (1981). Transforming genes of carcinomas and neuroblastomas introduced into mouse fibroblasts. *Nature* 290, 261-264.

Simpson, J.T., Workman, R.E., Zuzarte, P.C., David, M., Dursi, L.J., and Timp, W. (2017). Detecting DNA cytosine methylation using nanopore sequencing. *Nat Methods* 14, 407-410.

Sinn, E., Muller, W., Pattengale, P., Tepler, I., Wallace, R., and Leder, P. (1987). Coexpression of MMTV/*v*-Ha-ras and MMTV/*c*-myc genes in transgenic mice: synergistic action of oncogenes in vivo. *Cell* 49, 465-475.

Soneson, C., Love, M.I., and Robinson, M.D. (2015). Differential analyses for RNA-seq: transcript-level estimates improve gene-level inferences. *F1000Res* 4, 1521.

Stevenson, M., and Volsky, D.J. (1986). Activated *v*-myc and *v*-ras oncogenes do not transform normal human lymphocytes. *Mol Cell Biol* 6, 3410-3417.

Subramanian, A., Tamayo, P., Mootha, V.K., Mukherjee, S., Ebert, B.L., Gillette, M.A., Paulovich, A., Pomeroy, S.L., Golub, T.R., Lander, E.S., *et al.* (2005). Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc Natl Acad Sci U S A* 102, 15545-15550.

Takahashi, K., Tanabe, K., Ohnuki, M., Narita, M., Ichisaka, T., Tomoda, K., and Yamanaka, S. (2007). Induction of pluripotent stem cells from adult human fibroblasts by defined factors. *Cell* *131*, 861-872.

Taylor, A.M., Shih, J., Ha, G., Gao, G.F., Zhang, X., Berger, A.C., Schumacher, S.E., Wang, C., Hu, H., Liu, J., *et al.* (2018). Genomic and Functional Approaches to Understanding Cancer Aneuploidy. *Cancer Cell* *33*, 676-689 e673.

van der Maaten, L., and Hinton, G. (2008). Visualizing Data using t-SNE. *Journal of Machine Learning Research* *9*, 2579--2605.

Villanueva, A., Alsinet, C., Yanger, K., Hoshida, Y., Zong, Y., Toffanin, S., Rodriguez-Carunchio, L., Sole, M., Thung, S., Stanger, B.Z., *et al.* (2012). Notch signaling is activated in human hepatocellular carcinoma and induces tumor formation in mice. *Gastroenterology* *143*, 1660-1669 e1667.

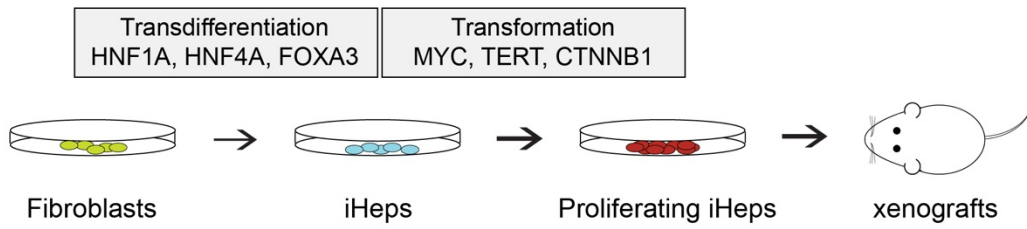
Vogelstein, B., Papadopoulos, N., Velculescu, V.E., Zhou, S., Diaz, L.A., Jr., and Kinzler, K.W. (2013). Cancer genome landscapes. *Science* *339*, 1546-1558.

Zerbino, D.R., Achuthan, P., Akanni, W., Amode, M.R., Barrell, D., Bhai, J., Billis, K., Cummins, C., Gall, A., Giron, C.G., *et al.* (2018). Ensembl 2018. *Nucleic Acids Res* *46*, D754-D761.

Zhu, B., Sun, L., Luo, W., Li, M., Coy, D.H., Yu, L., and Yu, W. (2017). Activated Notch signaling augments cell growth in hepatocellular carcinoma via up-regulating the nuclear receptor NR4A2. *Oncotarget* *8*, 23289-23302.

Figure 1.

A



B

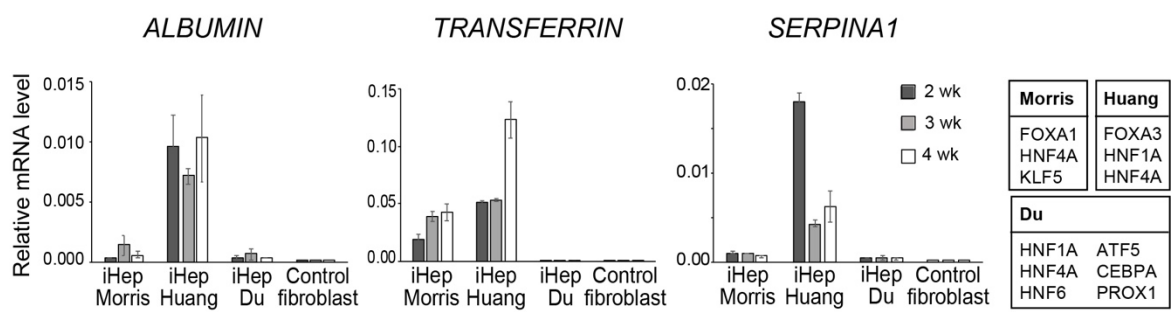


Figure 2.

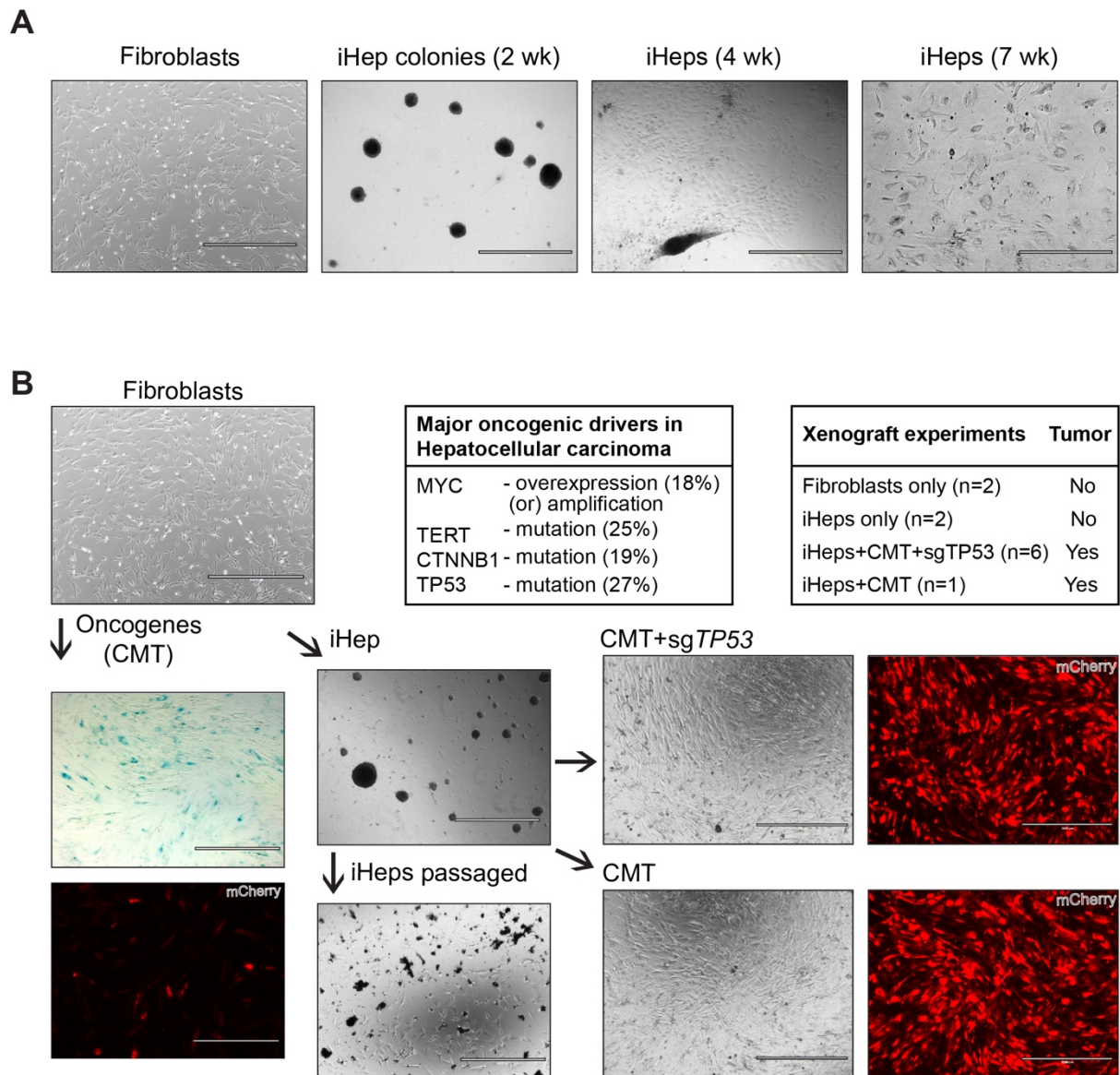


Figure 3.

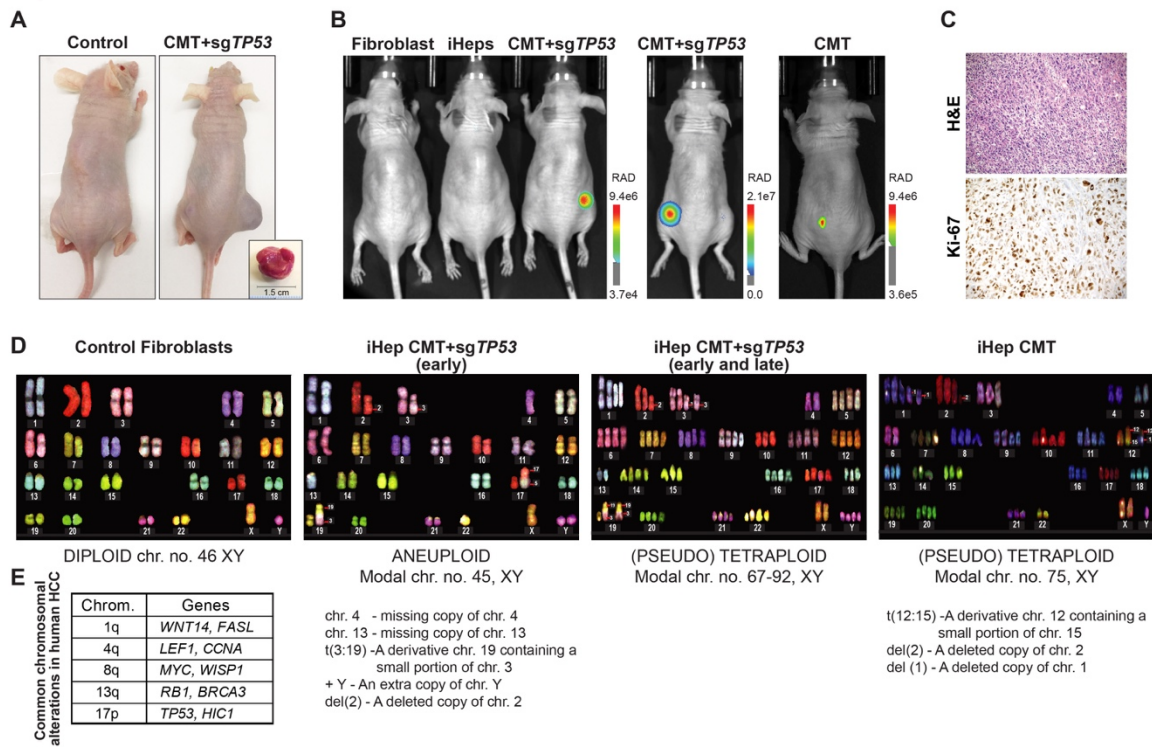


Figure 4.

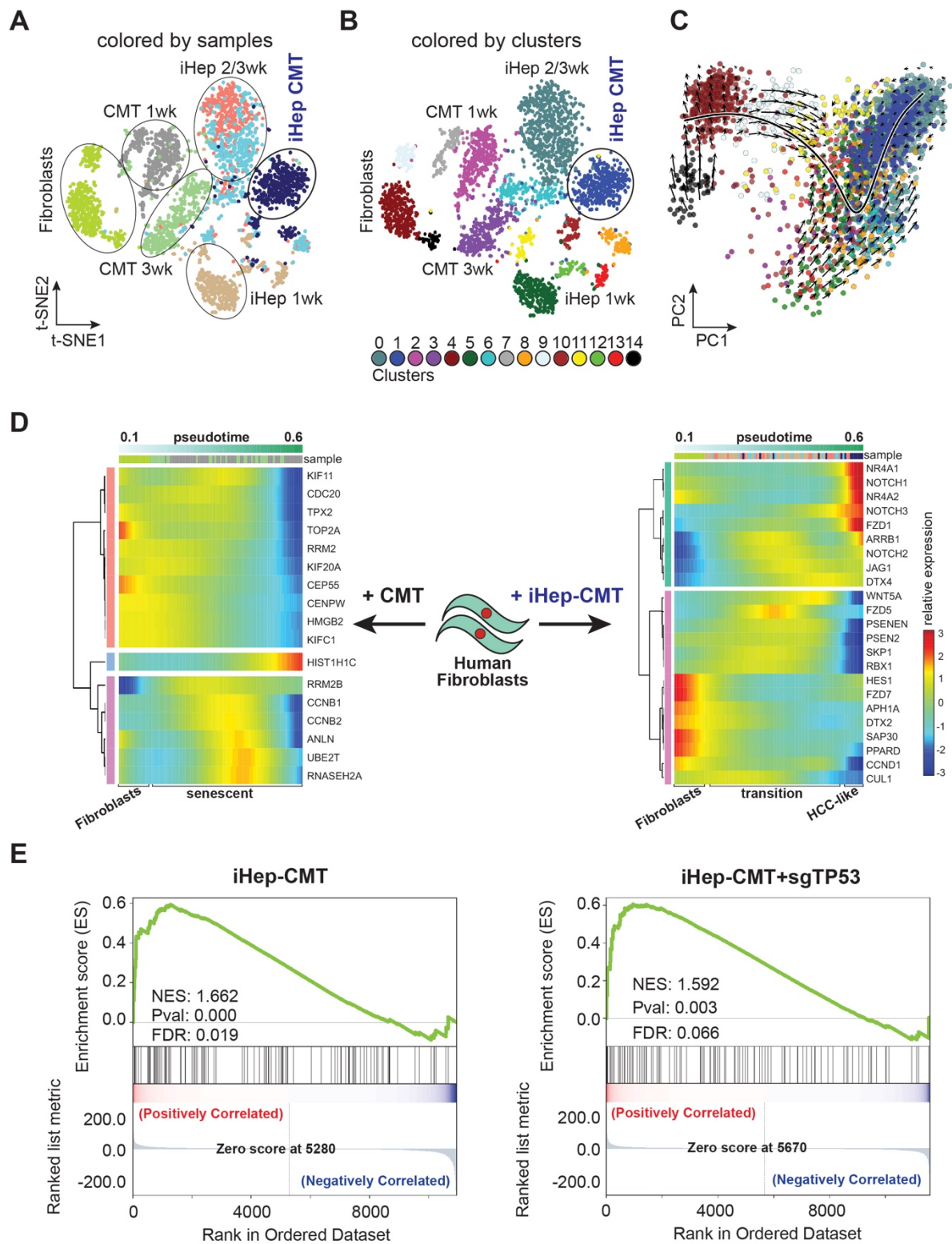
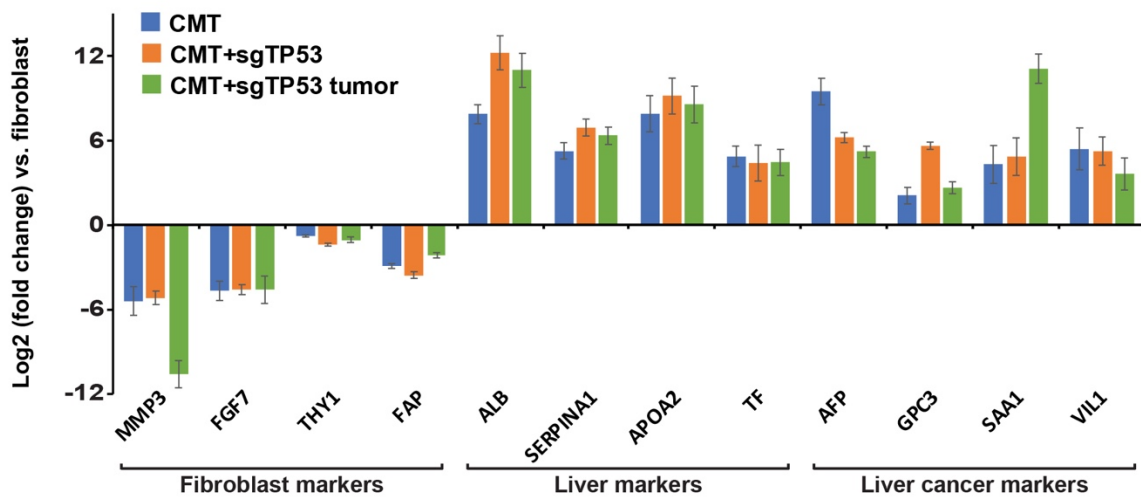
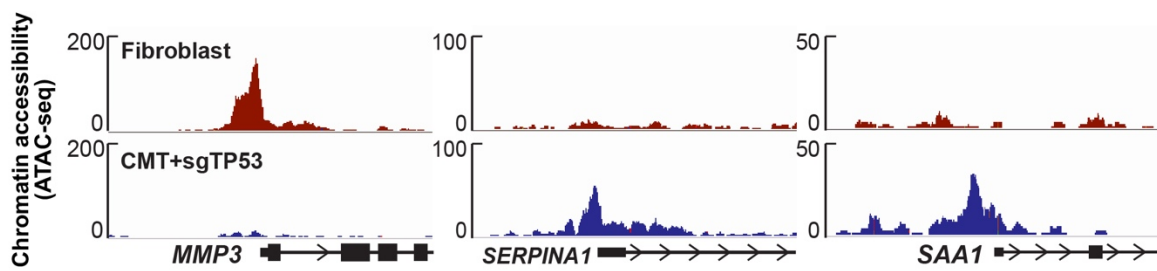


Figure 5.

A



B



C

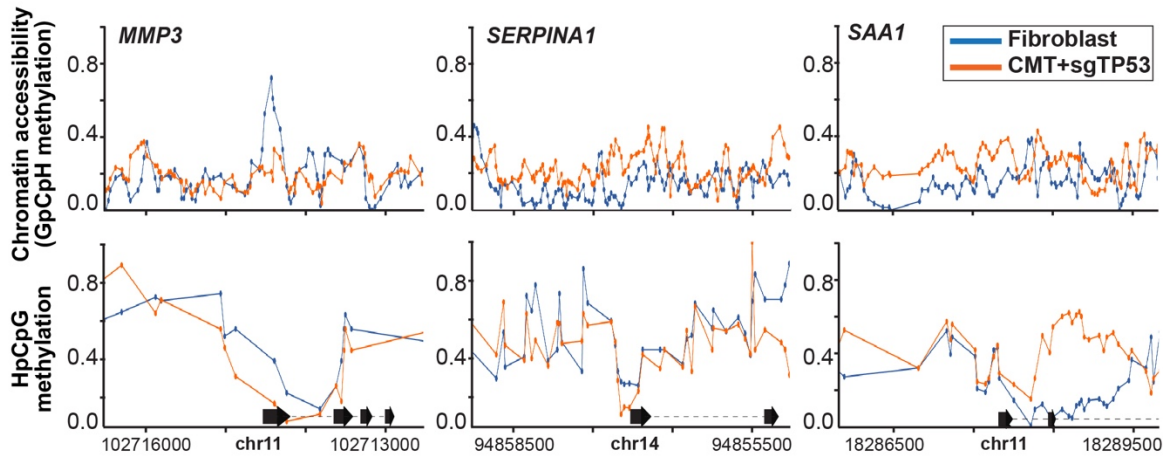
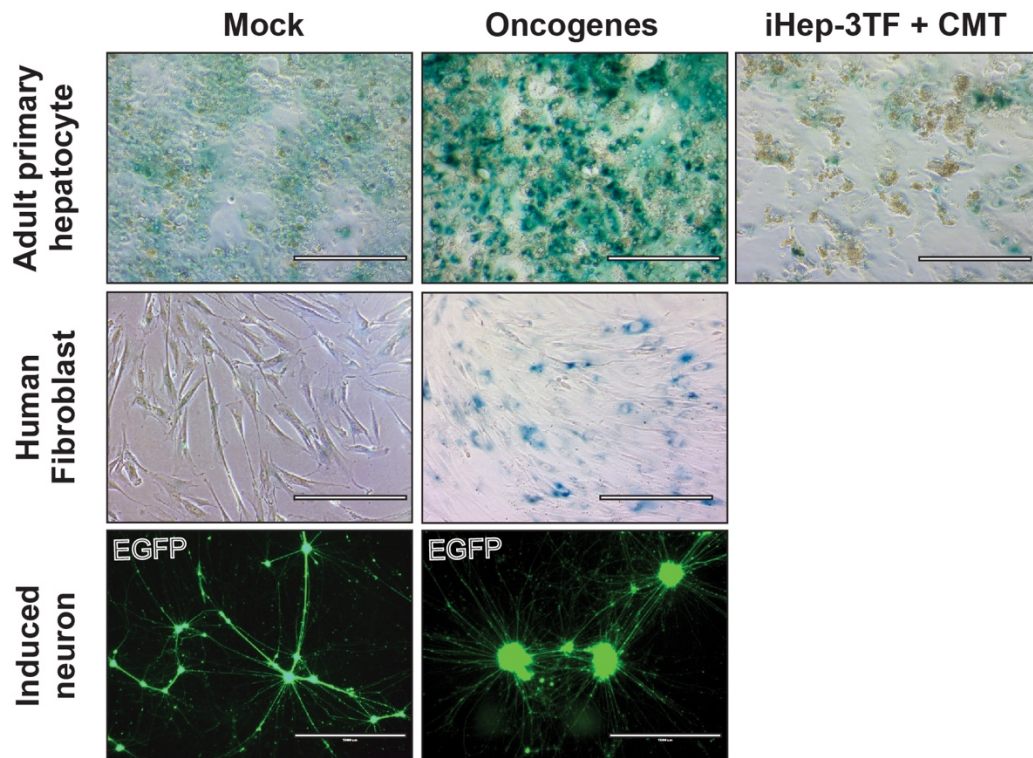


Figure 6.

A



B

