

# Pan-cancer analysis of clonality and the timing of systemic spread in paired primary tumors and metastases

Zheng Hu<sup>1,2,3</sup>, Zan Li<sup>4</sup>, Zhicheng Ma<sup>1,2,3</sup>, Christina Curtis<sup>1,2,3</sup> ‡

<sup>1</sup> Department of Medicine, Division of Oncology, Stanford University School of Medicine, Stanford, California, USA

<sup>2</sup> Department of Genetics, Stanford University School of Medicine, Stanford, California, USA

<sup>3</sup> Stanford Cancer Institute, Stanford University School of Medicine, Stanford, California, USA

<sup>4</sup> Life Science Research Center, Core Research Facilities, Southern University of Science and Technology, Shenzhen, Guangdong, China

‡ Correspondence: [cncurtis@stanford.edu](mailto:cncurtis@stanford.edu)

## Abstract

Metastasis is the primary cause of cancer-related deaths, but the natural history, clonal evolution and patterns of systemic spread are poorly understood. We analyzed exome sequencing data from 458 paired primary tumors (P) or metastasis (M) samples from 136 breast, colorectal and lung cancer patients, including both untreated (n=98) and treated (n=101) metastases. We find that treated metastases often harbored private driver gene mutations whereas untreated metastases did not, suggesting that treatment promotes clonal evolution. Polyclonal seeding was common in lymph node metastases (n=19/35, 54%; mostly untreated) and untreated distant metastases (n=20/70, 29%), but less frequent in treated metastases (n=9/90, 10%). The low number of metastasis-private clonal mutations is consistent with early metastatic seeding, which commonly occurred several years prior to diagnosis in breast (2.4 years, range 0–3.3), lung (3.6 years, range 2.8–3.7) and colorectal (4.1 years, range 3.1–4.6) cancers. Thus, this pan-cancer analysis reveals early systemic spread in three common cancer types. Further, these data suggest that the natural course of metastasis is selectively relaxed relative to early tumor development and that metastasis-private mutations are not drivers of cancer spread but are instead associated with drug resistance.

## Introduction

Metastasis remains poorly understood despite its critical clinical importance. For instance, metastases have been reported to originate from a single cell or clone in the primary tumor (monoclonal seeding)<sup>1-4</sup> or multiple clones (polyclonal seeding)<sup>5-7</sup>, but the prevalence of these patterns across distinct tumor types is unknown as is the impact of therapy and the timing of metastatic seeding<sup>8-10</sup>. While several recent studies have genomically characterized metastatic lesions<sup>11-13</sup> in the absence of the matched primary tumor, it is not feasible to disentangle the drivers of metastasis from those that are treatment associated since metastases are often sampled after treatment. However, comparisons of paired primary tumors and metastases have been far more limited due to the challenge in obtaining such samples<sup>5,8,14-18</sup>. As such, there has yet to be a systematic analysis of monoclonal versus polyclonal seeding, the chronology of systemic spread and the effect of therapy across cancers.

Here we analyzed whole-exome sequencing (WES) data from 458 paired primary tumor (P) and metastases (M) from 136 patients with colorectal, lung or breast cancers using a uniform bioinformatics pipeline. We assessed driver gene heterogeneity and evaluated the prevalence of monoclonal versus polyclonal seeding, revealing considerable variability between untreated and treated metastases across cancer types. Treatment was associated with high primary tumor versus metastasis (P/M) driver gene heterogeneity and monoclonal metastases. Metastatic seeding was estimated to occur two to four years prior to diagnosis of the primary tumor across three common cancer types, with breast cancers generally disseminating later and therefore closer to the time of detection relative to colorectal and lung cancers. Collectively, these observations suggest that systemic spread can begin early during tumor growth and that clonal architecture is remodeled by treatment, providing new insights into the clonal evolution of metastasis.

## Results

### The landscape of genomic alterations in paired primary tumors and metastases

We performed a literature review to identify cohorts with genomic sequencing data from matched normals, primary tumors (P) and metastases (M) from patients with three common cancer types, namely, colorectal<sup>16,17,19-22</sup>, lung<sup>23,24</sup> and breast<sup>23,25-29</sup> (**Table S1, Fig. S1**). All samples were processed within a uniform bioinformatics pipeline<sup>16,30</sup> to identify somatic single nucleotide variants (SSNVs), insertions/deletions (indels) and somatic copy number alterations (SCNAs) (**Methods**). Tumor purity/ploidy and cancer cell fraction (CCF) of SSNVs and indels (referred as SSNVs hereafter) were estimated in order to distinguish clonal (the upper bound of 95% confidence interval or CI of CCF  $\geq 1$ ) versus subclonal (the upper bound of 95% CI of CCF  $< 1$ ) SSNVs (**Methods**). Following quality

control assessment (**Methods**), we retained 458 tumor samples from 136 patients (colorectal cancer, n=39; lung cancer, n=30; breast cancer, n=67) for downstream analysis (**Table S1, Fig. S1**).

Overall, the mutational burden (SSNVs or SCNAs), tumor ploidy and SCNA frequency was highly concordant between P/M pairs (**Figs. S2-4, Table S2**), although differences between cancer types were noted. For instance, in breast cancer, the SCNA burden between P/M pairs was highly concordant but SSNV burden was only moderately concordant (**Fig. S2**), consistent with breast cancer being a copy number driven malignancy<sup>31</sup>. In all three cancer types, metastases exhibited a slight increase in the number of clonal SSNVs and fewer subclonal SSNVs (**Fig. S3**), consistent with an evolutionary bottleneck during metastasis. The mutational spectrum of M-private SSNVs (clonal or subclonal) between treated and untreated metastases was also highly concordant except that treated colorectal metastases were characterized by an enrichment of T>G transversions relative to untreated samples (**Fig. S5**). Indeed, all treated colorectal metastases (n=7) were biopsied after 5-fluorouracil (5-FU) chemotherapy in this cohort, which was recently shown to be associated with this mutational pattern<sup>32,33</sup>.

We next evaluated the enrichment of functional driver gene mutations in paired primary tumors and metastases. Three methods, namely PolyPhen-2<sup>34</sup>, FATHMM-XF<sup>35</sup> and CHASMplus<sup>36</sup>, were employed to assess the functionality (“driverness”) of nonsynonymous SSNVs in putative driver genes according to TCGA and COSMIC (**Methods, Table S3**). In total, 1085 functional driver SSNVs/indels were detected across these three cancer types (**Fig. 1a-b, Table S4**), in which 84%, 86% and 59% clonal drivers (including shared clonal, P-private clonal or M-private clonal) and 20%, 50% and 23% subclonal drivers (including shared subclonal, P subclonal/M clonal, P-private subclonal or M-private subclonal) were shared by P/M pairs for colorectal, lung and breast cancer, respectively (**Fig. 1c**). Amongst all driver mutations, M-private clonal and subclonal driver mutations were significantly enriched in breast cancer than colorectal and lung cancers (**Fig. 1c**). Gene ontology (GO) analysis of M-private driver genes revealed enrichment for chromatin binding, modification and organization genes (**Fig. S6, Table S5**), implicating chromatin regulators in metastatic progression<sup>37</sup>.

Amongst all non-silent clonal SSNVs in metastases, functional driver mutations were highly enriched on the trunk (P/M shared clonal) of the phylogenetic tree in both colorectal and breast cancers (**Fig. 1d, Methods**). However, this pattern was much weaker in lung cancer (**Fig. 1d**), presumably due to the large number of tobacco-associated non-silent clonal SSNVs (C>A mutations) induced early during lung cancer development (**Fig. S7**) as most of the lung cancer patients in this cohort (~90%) had a smoking history<sup>23,24</sup>. In

line with these results, the decreased ratio of nonsynonymous versus synonymous SSNVs (dN/dS)<sup>38</sup> in metastases (**Fig. S8**) suggests relaxed selective pressure relative to early cancer development in colorectal and breast cancers, but not lung cancer. Only 25%, 33% and 48% of colorectal, lung and breast cancer metastases, respectively, harbored one or more private clonal driver mutations (**Fig. 1e**) and these values were lower when restricted to untreated metastases (19%, 22% and 22%, respectively). Copy number analysis revealed a small number of putative driver genes that were more frequently amplified or deleted in metastases relative to paired primary tumors (increasing from P to M by 15%, **Fig. S9**). These include amplification of *RAC1* and deletions of *FAT1* and *ALB* in colorectal cancer, amplifications of *PLCG1* and *SALL4* and deletions of *NOTCH2*, *CDKN1B* in lung cancer and amplifications of *IL7R*, *NIPBL* and deletions of *NOTCH1*, *PTEN* in breast cancer (**Fig. S9**). Collectively, these data suggest that the genomic drivers required for invasion and metastasis often occur early in the primary tumor (**Fig. 1f**). Amongst treated metastases, the proportion of private-clonal drivers increased dramatically across all three cancer types with 71%, 75% and 53% in colorectal, lung and breast cancer, respectively (**Fig. 1e**). Amongst treated metastases, the proportion of private-clonal drivers increased dramatically across all three cancer types with 71%, 75% and 53% in colorectal, lung and breast cancer, respectively (**Fig. 1e**). This pattern was similarly evident in patients where both untreated and treated metastases were sampled (**Table S2**) where all (10/10) treated metastases harbored private functional driver mutation(s), but few (2/10) untreated lymph node metastases did (**Table S4**). Therefore, these data suggest that relapse after adjuvant therapy arises from a minor subclone in the primary tumor (**Fig. 1g**). In contrast, untreated metastases likely originate from the major clone in the primary tumor (**Fig. 1e**). Hence, treatment confers a stringent selective pressure and promotes clonal evolution of the metastasis.

## Patterns of metastatic seeding in lymph node and distant metastases

In order to infer the clonality of individual metastases (**Fig. 2a**), we compared the CCFs of SSNVs in each P/M pair and the number of M-private clonal SSNVs, P-private clonal SSNVs and P/M shared subclonal SSNVs was denoted as  $L_m$ ,  $L_p$  and  $W_s$ , respectively (**Fig. 2b**). We used the Jaccard similarity index (JSI) where  $JSI = W_s / (L_m + L_p + W_s)$  to quantify mutational similarity between P/M pairs<sup>39</sup> (**Methods**). Polyclonal seeding is expected to result in a higher JSI than monoclonal seeding due to the higher proportion of shared subclonal SSNVs (higher  $W_s$ ) and the presence of fewer M or P-private clonal SSNVs (lower  $L_m$  and  $L_p$ ) (**Fig. 2b**). These patterns were verified by simulation studies using an established agent-based model of spatial tumor progression<sup>16,30</sup> (**Figs. S10-S11, Methods**). By analyzing data from virtual tumors simulated under varied parameters (**Methods**), we found that a JSI value of 0.3 maximizes the classification accuracy (91.1%) in distinguishing monoclonal versus polyclonal seeding (**Fig. 2c**). Hence, this cutoff was applied to the patient genomic data (**Fig. 2c**). Most metastases exhibited patterns

consistent with monoclonal seeding ( $n=151$ , 76% of metastases; median JSI=0.075, interquartile range, IQR=0.021–0.138), whereas polyclonal seeding was less frequent ( $n=48$ , 24% of metastases; median JSI=0.523, IQR=0.469–0.800) (**Figs. 2c**).

As expected, monoclonal metastases ( $n=151$ ) exhibited significantly higher  $L_m$  and  $L_p$  values than polyclonal metastases ( $n=48$ ) ( $P=6.2e-16$  and  $P=2.1e-09$  for  $L_m$  and  $L_p$ , respectively, two-sided Wilcoxon Rank Sum Test) and significantly lower  $W_s$  values ( $P=2.1e-12$ , two-sided Wilcoxon Rank Sum Test) (**Fig. 2d**). Metastases of monoclonal origin also harbored significantly more SCNAs relative to paired primary tumors than polyclonal metastases ( $P=1.9e-08$ , two-sided Wilcoxon Rank Sum Test; **Fig. 2e**). Indeed,  $L_m$  is highly correlated with the number of P-to-M altered SCNAs (Pearson's  $R=0.52$ ,  $P=5.0e-15$ ; **Fig. 2f**), indicating that both SSNVs and SCNAs reflect the clonality of metastases. Polyclonal seeding was more prevalent in axillary lymph node metastases (19/35 or 54%) relative to distant metastases (29/164 or 18%) ( $P=1.8e-05$ , two-sided Fisher's exact test; **Figs. 2g and S12a**), potentially reflecting greater lymphatic spread of disseminated cells to the lymph nodes. Amongst distant metastases, polyclonal seeding was more prevalent in untreated metastases (20/70 or 29%) than treated metastases (9/90 or 10%) ( $P=0.002$ , two-sided Fisher's exact test; **Fig. 2g**), presumably because treatment selects for resistant subclones that dominate the relapse resulting in monoclonal metastases (**Fig. 2h**). The higher P/M driver gene heterogeneity observed in treated versus untreated metastases (**Fig. 1e**) is consistent with this scenario. The prevalence of polyclonal seeding differed across metastatic sites (lymph node, liver, brain and lung), with brain and lung more commonly exhibiting monoclonal seeding (**Fig. S12b**); these two sites were more commonly biopsied after treatment. We verified the JSI-based classification of monoclonal versus polyclonal seeding by phylogenetic analysis of patients with multi-region sequencing (MRS) data of the primary tumor and metastasis ( $n=13$  patients; **Figs. 3 and S13**). Monoclonal seeding was associated with a monophyletic tree structure (metastatic samples make up a single phylogenetic clade), whereas polyclonal seeding was associated with a polyphyletic structure (metastatic samples make up multiple phylogenetic clades) (**Figs. 3 and S13**).

## Chronology of metastatic seeding

Previously, we described a computational framework (SCIMET) to estimate the timing of metastatic seeding relative to primary tumor size based on multi-region sequencing (MRS) of P/M pairs<sup>16</sup>. Application of this approach to colorectal cancer yielded quantitative evidence for early systemic spread, well before the primary tumor was clinically detectable. Since MRS data was not available for the vast majority of patients in this cohort, we developed a new computational method that leverages exome sequencing data from a single biopsy to time metastatic seeding (**Figs. 4a and S14, Supplementary Note**). The



time (in years) from metastatic seeding to diagnosis of the primary tumor ( $t_s$ ) can be approximated by:

$$t_s \approx (1 - \frac{L_m}{L_p} \alpha) \times T \quad \text{Eq.(1)}$$

where  $L_m$  and  $L_p$  correspond to the number of M-private clonal SSNVs and P-private clonal SSNVs, respectively;  $T$  is the primary tumor expansion age (time from emergence of carcinoma founder cell to diagnosis);  $\alpha = t_p/T$  where  $t_p$  is the time from emergence of carcinoma founder cell to the most recent common ancestor in the primary tumor sample (pMRCA, **Figs. 4a** and **S14, Supplementary Note**). The time fraction  $\alpha$  is expected to be small because bulk sequencing only detects relatively high frequency mutations that occur early during tumor growth or are strongly selected for<sup>40-42</sup>. We applied our established agent-based model of spatial tumor growth<sup>30</sup> to simulate a large set of virtual tumors (n=1000, each  $\sim 10^9$  cells) with varying growth rates (**Methods**). *In silico* sequencing of a single biopsy (each  $\sim 10^6$  cells, mean depth=100X) from the virtual tumors (n=1000) yields an estimate of  $\tilde{\alpha}=0.13\pm 0.0028$  (**Fig. S14**), confirming the observation that bulk sequencing typically only detects high-frequency mutations that occur early during tumor growth. Here we assume a model of stringent selection (selection coefficient,  $s=0.1$ ) during growth of the primary tumor since most primary tumors in this cohort (57/65 or 88% evaluable tumors) exhibited variant allelic frequencies (VAF) that were not consistent with neutral evolution<sup>43</sup> (**Fig. S15; Methods**).

We utilized a Gompertzian model of tumor growth<sup>44</sup>, to estimate the tumor expansion age ( $T$ ) for each of the three cancer types (**Supplementary Note**) where tumor size and doubling time (DT) at diagnosis were obtained from literature review (**Table S6**). This yields estimates of average tumor expansion age of  $\tilde{T}=5.2$  (IQR, 4.3–7.7), 4.3 (IQR, 2.7–4.4) and 4.6 (IQR, 3.2–6.6) years for colorectal, lung and breast cancer, respectively (**Fig. 4b** and **Table S7**). Chronological estimates of seeding time relative to diagnosis of the primary tumor ( $\tilde{t}_s$ ) can be computed by Eq.(1) as follows: 4.1 years (IQR, 3.2–4.6), 3.6 years (IQR, 2.8–3.7) and 2.7 years (IQR, 1.1–3.5) for colorectal, lung and breast cancers, respectively (**Fig. 4c** and **Table S7**). The estimated timing of metastasis here ( $\tilde{t}_s$ ) agreed with our previous estimates (using the colorectal cancer cohort) of primary tumor size at time of metastatic seeding<sup>16</sup> ( $R= -0.58$ ,  $P=0.009$ , **Fig. S16**; note the negative correlation since this study estimates backward time and the previous study estimates forward time). Of note, while  $\tilde{t}_s < 0$  may indicate metastatic seeding after diagnosis/resection of the primary tumor, large  $L_m$  values can lead to  $\tilde{t}_s < 0$  (see Eq.(1)) even when the metastasis was seeded before diagnosis of the primary tumor. To mitigate this uncertainty, samples with estimated seeding times later than the actual time of diagnosis of metastasis were excluded (n=12 for breast, 1 for colorectal and 1 for lung cancer, respectively) (**Supplementary Note**). We find that  $\tilde{t}_s < 0$  was more common in breast cancer and more generally breast cancers disseminated closer to the time of detection (later) compared to colorectal and lung cancers (**Fig. 4c**). This may be because screening mammography

detects relatively small primary breast tumors (<2 cm)<sup>45</sup>. However, even after normalization to primary tumor age (namely  $t_s/T$ ), which depends on tumor size and the underlying growth parameters (**Supplementary Note**), breast cancer was found to disseminate later than colorectal and lung cancers (**Fig. S17**). Most breast cancer metastases (83%) in this cohort were biopsied after adjuvant therapy (**Fig. S1**), whereas this fraction is fewer in colorectal (13%) and lung (20%) cancer metastases and breast cancers harbored more private driver mutations than colorectal and lung cancers (**Figs. 1a-c**). Thus, the genomic complexity of metastatic relapses in breast cancer relative to unpaired early-stage primary tumors<sup>12</sup> at least in part reflects the selective effect of treatment on the genome rather than the drivers of metastatic spread. Of note, HER2-positive breast cancers tended to disseminate earlier than HER2-negative breast cancers (**Fig. S18**) consistent with this subgroup having the highest risk of distant metastasis before the routine use of the HER2-targeted therapy, trastuzumab, which has revolutionized the treatment of this disease<sup>46</sup>.

As expected, metachronous metastases were often seeded later than synchronous metastases (median  $t_s=3.8$  vs 3.0,  $P=5.6e-05$ , two-sided Wilcoxon Rank-Sum Test; **Fig. 4d**). In fact,  $t_s$  was highly correlated with the clinical time span from diagnosis of primary tumor to metastasis (**Fig. 4e**), indicating that metastases that manifest late clinically were seeded later. Since primary tumor size at diagnosis is an important predictor of a patient's prognosis (time to metastatic relapse) (**Fig. S19a**), we suspect that metastases in patients with larger primary tumor size at initial diagnosis were seeded earlier (namely larger  $t_s$ ). Indeed,  $t_s$  is positively associated with the primary tumor size at diagnosis ( $R=0.32$ ,  $P=0.00024$ ; **Fig. S19b**). These results corroborate our estimates of metastatic timing. According to *Eq.(1)*, a larger number of M-private clonal mutations (larger  $L_m$ ) indicates later dissemination. Supporting this theory, metachronous metastases showed significantly larger  $L_m$  than synchronous metastases (metachronous: median  $L_m=24$ , IQR=16–40; synchronous: median  $L_m=11$ , IQR=6–32;  $P=6.5e-4$ , two-sided Wilcoxon Rank-Sum Test; **Fig. S20a**). This pattern held for SCNAs where metachronous metastases showed significantly more SCNAs relative to the primary tumor as compared to synchronous metastases (**Fig. S20b**). Since metachronous metastases were generally seeded later than synchronous metastases (**Fig. 4d**), this is consistent with the higher degree of genomic divergence with primary tumor in late seeded metastases<sup>18</sup>. Collectively, these data indicate that systemic spread can occur several years prior to diagnosis of the primary tumor but with variability across cancer histologies and subtypes.

## Discussion

We performed a systematic analysis of exome sequencing data in paired primary tumors and metastases across three common cancers: colorectal, lung and breast and find that polyclonal seeding is common in lymph node metastases (19/35, 54%; most untreated)

and untreated distant metastases (20/70, 29%), but rare (9/94, 10%) in metastases sampled after adjuvant therapy (**Fig. 2g**). Consistent with these results, treated metastases were strongly enriched for functional driver mutations as compared to untreated metastases (**Fig. 1e**). This finding indicates that driver gene heterogeneity is minimal between untreated metastases and primary tumors (**Fig. 1e**). Comparisons of paired primary tumors and distant metastases indicates that systemic spread can occur rapidly following malignant transformation, often several years prior to diagnosis of the primary tumor across three major types (**Fig. 4c**). These results are consistent with other reports of early seeding based on animal models and disseminated tumor cells <sup>9,47,48</sup>.

Our analyses on driver gene heterogeneity, clonality and the timing of metastases provide important insights into the clonal dynamics of metastatic progression. First, in the absence of treatment, metastases often arise from the major clone in the primary tumor and lack metastasis-specific driver mutations (**Fig. 1f**). Consistent with these observations, a recent pan-cancer study demonstrated that driver gene heterogeneity is also minimal amongst multiple untreated metastases <sup>49</sup>. Moreover, the prevalence of polyclonal seeding in untreated lymph node and distant metastases indicates multiple cell subpopulations in primary tumor have acquired the metastatic competence. Half of all metastases (51%) studied here were biopsied after treatment, and these commonly exhibited monoclonal seeding accompanied by private driver mutations. As such, polyclonal seeding may be relatively common, but the ultimate pattern of clonality in the metastatic lesion is influenced by treatment.

Second, our quantitative framework demonstrates that systemic spread typically begins 2-4 years prior to the diagnosis of primary tumor (**Fig. 4c**). These data suggest that in some patients, metastatic seeding can happen very early especially for synchronously diagnosed metastases (**Figs. 4e, 5a**). Metachronous distant metastases following treatment occurred relatively later than synchronous distant metastases and harbored more genomic variations and driver mutations (**Figs. 1e and S20**). These data suggest that treatment remodels the clonal evolution of metastasis by selecting disseminated cells with drug resistant mutations (**Fig. 5a-b**). As such metastasis-specific mutations are unlikely to be the drivers of metastasis, but instead are associated with drug resistance (**Fig. 5b**). This interpretation is of clinical relevance and helps to clarify the observation that metastatic relapses are more genomically complex than unpaired early-stage primary breast tumors <sup>12</sup>.

We also observe that many breast cancer relapses are diagnosed well after (>5 years) initial diagnosis, although they were seeded several years prior to primary tumor diagnosis (**Fig. 4e**). This may reflect periods of quiescence or dormancy of disseminated tumor cells (DTCs) <sup>50,51</sup> (**Fig. 5a**). Indeed, ongoing research is focused on the development of dormancy-targeted therapies to prevent metastatic relapse <sup>52,53</sup>. Here we analyzed



clinically detectable metastases and it is impossible to know how many DTCs and micrometastases were eliminated as a result of adjuvant treatment (**Fig. 5b**) or immune surveillance. Indeed, an important area of ongoing research is to elucidate how the immune system can be harnessed and whether immunotherapy could be administered in the adjuvant setting to prevent metastatic progression.

## Acknowledgments

We thank Hang Xu, Katherine McNamara, Eran Kotler, Jennifer Caswell-Jin and other members of Curtis laboratory for valuable discussions. We thank Jiguang Wang and Quanhua Mu for providing the scripts for the ternary plot. C.C. is supported by the National Institutes of Health through the NIH Director's Pioneer Award (DP1-CA238296), the American Association for Cancer Research (AACR) and the Emerson Collective. Z.H is supported by an Innovative Genomics Initiative (IGI) Postdoctoral Fellowship.

## Author contributions

Z.H and C.C conceived and designed the study. Z.H performed all computational analyses. Z.L reviewed the published studies, extracted and analyzed the clinical data. Z.M processed the in-house clinical samples and generated the genomic data. Z.H and C.C wrote the manuscript, which was reviewed by all authors.

## Methods

### Whole-exome sequencing (WES) of paired primary tumors and metastases

We performed a comprehensive review on the published studies through surveying the PubMed database (<https://www.ncbi.nlm.nih.gov/pubmed/>), in which whole-exome sequencing (WES) was performed for matched normal tissues, primary tumors (P) and metastases (M) in the same patients. We focused on colorectal, lung and breast cancers given the availability of large patient data in these three cancer types. In total, the raw sequencing reads data for 586 tumor samples from 181 patients in 13 published studies were accessed and retrieved (**Table S1**). We also generated multi-region sequencing (MRS) data for two colorectal cancer patients (mCRCTB1 and mCRCTB7) with liver metastases for whom multi-region sequencing (MRS) data (n=5-7 sites for each of P and M; 24 tumor samples in total). Here tumor tissues with cellularity >60% were selected for DNA isolation using the QIAamp DNA FFPE Tissue Kit (Qiagen) and libraries were generated using the Agilent SureSelect Human All Exon kit for sequencing on the Illumina HiSeq 2500. In total, the WES data in 610 tumor samples from 183 metastatic cancer patients including 54 colorectal cancers (215 tumor samples), 35 lung cancer (87 tumor samples) and 94 breast cancers (308 tumor samples) were analyzed in this study. Clinical information was retrieved from the original studies, including patient age at initial diagnosis, time span from initial diagnosis of primary tumor to diagnosis of metastasis, treated

information, cancer subtypes, *etc.* (**Table S2**). We define synchronous metastases if the time span between diagnosis of primary tumor and metastasis is within 3 months and metachronous metastases if the time span is  $\geq 3$  months.

An established bioinformatics pipeline was used to detect somatic single nucleotide variations (SSNVs), small insertions/deletions (indels) and somatic copy number alterations (SCNAs), estimate tumor purity/ploidy and estimate the cancer cell fraction (CCF) for each SSNVs/indels in corresponding samples <sup>16,30</sup>. In particular, paired sequencing reads were aligned to human reference genome (NCBI build hg19) with BWA (v.0.7.10) <sup>54</sup>. Duplicate reads were marked with Picard Tools (v.1.111). Aligned reads were further processed with GATK 3.4.0 for local re-alignment around insertions and deletions and base quality recalibration.

### SSNVs and indel calling

SSNVs were called by MuTect (v.1.1.7) <sup>55</sup> for each tumor/normal pair. SSNVs failing MuTect's internal filters, having fewer than 10 total reads or 3 variant reads in the tumor sample, fewer than 10 total reads in the normal sample, or mapping to paralogous genomic regions were removed. Additional VarScan (v.2.3.9) <sup>56</sup> filters were applied to remove SSNVs with low average variant base qualities, low average mapping qualities among variant supporting reads, strand bias among variant supporting reads and high average mismatch base quality sums among variant supporting reads, either within each tumor sample or across all tumor samples from the same patient. The maximal observed variant allele frequencies (VAF) across all samples from each patient were calculated based on raw output files from MuTect. SSNVs with maximal observed VAFs lower than 0.05 were removed. For FFPE specimens, additional filters were applied to exclude possible artifactual SSNVs. Specifically, artifacts among C>T/G>A SSNVs with bias in read pair orientation were filtered in each individual FFPE sample, similar to the approach of Costello *et al* <sup>57</sup>. We also sought to exploit the multi-sample information in the same patients to retrieve read counts for SSNVs. To obtain the depth and VAF information across all samples from the same patient, for each SSNV and in each tumor sample that an SSNV was not originally called in, the total reads and variant supporting reads were counted using the *mpileup* command in SAMtools (v.1.2) <sup>58</sup>. Only reads with mapping quality  $\geq 40$  and base quality at the SSNV locus  $\geq 20$  were counted and used to calculate the VAF for that SSNV. Small insertions/deletions (indels) were called with Strelka (v.1.0.14) <sup>59</sup>. SSNVs and indels were annotated with ANNOVAR (v.20150617) <sup>60</sup> and those in protein coding regions were retained for downstream analyses.

### Copy number analysis

Copy number analysis was performed using TitanCNA (v.1.5.7) <sup>61</sup>. Briefly, TitanCNA uses depth ratio and B-allele frequency information to estimate allele-specific absolute copy

numbers with a hidden Markov model, and estimates tumor purity and clonal frequencies. Only autosomes were used in copy number analysis. First, for each patient, germline heterozygous SNP at dbSNP 138 loci were identified using SAMtools and SnpEff (v.3.6) in the normal sample. HMMcopy (v.0.99.0) <sup>62</sup> was used to generate read counts for 1000bp bins across the genome for all tumor samples. TitanCNA was used to calculate allelic ratios at the germline heterozygous SNP loci in the tumor sample and depth ratios between the tumor sample and the normal sample in bins containing those SNP loci. Only SNP loci within WES covered regions were then used to estimate allele-specific absolute copy number profiles. TitanCNA was run with different numbers of subclones (n=1-3). One run was chosen for each tumor sample based on visual inspection of fitted results, with preference given to the results with a single subclone unless results with multiple subclones had visibly better fit to the data. Results from tumor samples from the same patient were inspected together to ensure consistency. Overall ploidy and purity for each tumor sample was calculated from the TitanCNA results.

Differentially altered SCNAs in the metastasis relative to paired primary tumor (P-to-M) were identified if following three criteria were satisfied simultaneously: 1) absolute copy number in the metastasis was larger than 2.8 or less than 1.2; 2) copy number relative to median ploidy in the metastasis was larger than 0.8 or less than -0.8; 3) changes relative to the primary tumor in both absolute copy number and relative copy number were larger than 0.8 or less than -0.8.

### Cancer cell fraction (CCF) estimates and identification of clonal and subclonal mutations

The CCFs and their variation (95% confidence interval or 95% CI) for each SSNVs/indels in the corresponding samples were estimated with CHAT (v 1.0) <sup>63</sup>. CHAT includes a function to estimate the CCF of each SSNVs by adjusting its variant allele frequency (VAF) based on local allele-specific copy numbers at the SSNV locus. SSNV frequencies and copy number profiles estimated from previous steps were used to calculate the CCFs for all SSNVs in autosomes. The CCFs were also adjusted for tumor purity using the estimates by TitanCNA. In brief, for an SSNV residing in a genomic segment with a total copy number of  $CN_t$ , minor allele copy number of  $CN_b$  and cellular prevalence  $P_{CNA}$  of the CNA in the tumor content, the estimated CCF of the SSNV is:

$$CCF = \begin{cases} CN_c \times \frac{VAF}{p'} - P_{CNA} \times (CN_t - CN_b - 1) & \text{Early Major} \\ CN_c \times \frac{VAF}{p'} - P_{CNA} \times (CN_b - 1) & \text{Early Minor} \\ CN_c \times \frac{VAF}{p'} & \text{Late/Independent} \end{cases} \quad \text{Eq. (2)}$$

where  $CN_c = CN_t \times P_{CNA} + 2 \times (1 - P_{CNA})$  and the effective purity  $p' = \frac{CN_t \times p}{CN_t \times p + 2 \times (1 - p)}$  ( $p$  is estimated tumor purity) and VAF is the observed variant allele frequency. The temporal

ordering and background composition of SSNVs and SCNAs was inferred by comparing the conditional probabilities of the observed number of mutant reads out of total reads, under each scenario and CNA configuration ( $CN_t$ ,  $CN_b$ ,  $P_{CNA}$ ) as follows: *Early Major* or *Minor*: SSNV in the major or minor allele occurred before the CNA; *Late*: SSNV occurred after the CNA; *Independent*: the SSNV and CNA occurred in independent lineages<sup>63</sup>.

To distinguish clonal and subclonal SSNVs/indels in each sample, we employ the following criterion: clonal – 95% CI overlaps with 1; subclonal – the upper bound of 95% CI is smaller than 1, as previously used<sup>64</sup>. The CCFs of SSNV/indels for each P/M sample pair were visualized using the scatter plot and manually checked in order to identify problematic samples. In particular, for each P/M pair, a cluster of SSNV/indels centered around CCF=1 is expected which represent truncal (P-M shared clonal) mutations that occurred prior to malignant transformation of the founding cell in the primary tumor. The patients (n=5) with none of or very few (<10) trunk SSNVs/indels were excluded as which implies independent (non-clonal) origin for the primary tumor and metastasis. Furthermore, patients (n=42) with a diffusely distributed cluster for truncal SSNVs/indels were also excluded since this is likely caused by low tumor purity or low sequencing quality. After these filtering steps, 458 tumor samples from 136 metastatic cancer patients including 39 colorectal cancers (181 tumor samples), 30 lung cancer (75 tumor samples) and 67 breast cancers (202 tumor samples) were retained for downstream analysis in this study.

### Jaccard similarity index

The number of M-private clonal, P-private clonal and P-M shared subclonal SSNVs for each P/M pair was denoted as  $L_m$ ,  $L_p$  and  $W_s$  respectively. For two sets, the Jaccard similarity index (JSI) is defined for the intersection divided by the union of these two sets. Thus, the JSI for a P/M pair can be defined as:

$$JSI = \frac{W_s}{L_p + L_m + W_s} \quad Eq. (3)$$

For multi-region sequencing data,  $L_m$ ,  $L_p$  and  $W_s$  was counted by pairwise comparison of each sample pair from the P and M. The mean  $L_m$ ,  $L_p$  and  $W_s$  was used to compute the JSI by Eq.(3).

### Functional assessment of non-silent somatic mutations

To identify functional driver gene mutations, three commonly used computational methods, PolyPhen-2<sup>34</sup> (<http://genetics.bwh.harvard.edu/pph2/>), FATHMM-XF<sup>35</sup> (<http://fathmm.biocompute.org.uk/fathmm-xf/>) and CHASMplus<sup>36</sup> (<https://karchinlab.github.io/CHASMplus/>), were utilized to perform the function (“driverness”) assessment on the nonsynonymous SSNVs amongst putative cancer genes derived from TCGA pan-cancer<sup>65</sup> and COSMIC (Release v87, Nov. 13, 2018).

Stopgain/splicing point mutations and indels on putative cancer genes are classified as functional drivers automatically.

Putative cancer genes were curated by merging all TCGA pan-cancer drivers (n=299)<sup>65</sup> and additional cancer type-specific drivers annotated by COSMIC Cancer Gene Census (<https://cancer.sanger.ac.uk/cosmic>; n=47, 40 and 9 for colorectal, lung and breast cancers, respectively). For PolyPhen-2, a SSNV is considered as “functional” when the functional report (“pph2\_class”) is “deleterious”. For FATHMM-XF, a SSNV is considered as “functional” when the functional report (“Warning”) is “pathogenic”. For CHASMplus, a SSNV is considered as “functional” when the FDR < 0.05. In this study, the SSNVs, predicted to be functional by any of these three methods, were considered as functional mutations. Metascape<sup>66</sup> (<http://metascape.org>) was used to perform gene ontology (GO) analysis of functional driver genes.

### Driver enrichment analysis

Clonal non-silent SSNVs/indels in a metastatic lesion can be considered truncal clonal (or P-M shared clonal) or M-private clonal where the number is denoted  $L_s\_total$  and  $L_m\_total$ , respectively. Meanwhile, the functional driver SSNVs/indels in a metastasis are denoted  $L_s\_driver$  and  $L_m\_driver$ , respectively. The ratios,  $L_s\_total/L_m\_total$  and  $L_s\_driver/L_m\_driver$ , can be evaluated for functional enrichment of drivers on the truncal or M-private branch of the corresponding phylogenetic tree. Since  $L_s\_driver$  and  $L_m\_driver$  are small values ( $L_m\_driver \sim 0$  for many metastases), they lead to high variation in the  $L_s\_driver/L_m\_driver$  ratio. A down-sampling (bootstrapping) step (50% of the patients each time) was performed in which sampled patient data were merged to derive the  $L_s\_total/L_m\_total$  and  $L_s\_driver/L_m\_driver$  ratios. 100 repeated down-samplings were performed for each of the three cancer types to derive statistical measures.

### Mutational signatures, dN/dS and test of neutrality

MuSiCa<sup>67</sup> (<http://bioinfo.ciberehd.org:3838/MuSiCa/>) was used to extract mutation signatures based on non-negative matrix factorization<sup>68</sup> for P/M shared clonal (truncal) SSNVs, M-private clonal SSNVs and M-private subclonal SSNVs respectively, in each of the three cancer types. dndscv<sup>38</sup> (<https://github.com/im3sanger/dndscv>) was used to compute the ratio of nonsynonymous and synonymous SSNVs (dN/dS) for missense and nonsense mutations, respectively and for P/M shared clonal (trunk) SSNVs, M-private clonal SSNVs and M-private subclonal SSNVs, respectively, in each of the three cancer types. We evaluated whether a tumor follows neutral evolution or under strong selection during the growth by analyzing the variant frequency distribution (VAF) of subclonal SSNVs. Under neutral evolution<sup>43</sup>, the number of subclonal SSNVs with VAF larger than  $f$  in a tumor cell population follows a power-law distribution:  $m(f) \sim 1/f$ . The adjusted VAFs (equivalent to CCFs/2) for subclonal SSNVs (in the range of 0.1–0.3) were used here and only tumors harboring at least 20 subclonal SSNVs in this range were analyzed (n=65).



primary tumors and 79 metastases). By fitting this model and using a threshold of  $R^2=0.98$ , the mode of evolution (neutral or selection) can be inferred (**Fig. S15**).

### Phylogenetic tree reconstruction

We ran PHYLIP<sup>69</sup> via an online version (<http://www.trex.uqam.ca/index.php?action=phylip&apP=dnaps>) and applied the Maximum Parsimony method to reconstruct the phylogeny of multiple specimens from individual patients based on the presence or absence of SSNVs/indels. The SSNVs/indels residing a region with different loss-of-heterozygosity (LOH) status between paired primary tumor and metastasis were filtered, since which may lead to erroneous presence or absence of SSNVs/indels in paired P and M. When multiple maximum parsimony trees were reported, we chose the top ranked solution. FigTree (<http://tree.bio.ed.ac.uk/software/figtree/>) was employed to visualize the reconstructed trees.

### Spatial agent-based modeling of metastatic progression

We employed our previously established three-dimensional agent-based tumor evolution framework<sup>30</sup> to model tumor growth, mutation accumulation and metastatic dissemination after malignant transformation. Pre-malignant clonal expansions prior to transformation do not alter the genetic heterogeneity within a tumor thus were not modeled and we assume that dissemination occurs after malignant transformation of the founding carcinoma cell. In this model, spatial tumor growth is simulated via the expansion of deme subpopulations (composed of ~5k cells with diploid genome), mimicking the glandular structures often found in epithelial tumors and metastases and consistent with the number of cells found in individual colorectal cancer glands (~2,000-10,000 cells). The deme subpopulations expand within a defined 3D cubic lattice (Moore neighborhood, 26 neighbors), via peripheral growth while cells within each deme are well-mixed without spatial constraints and grow via a random birth-and-death process (division probability  $b$  and death probability  $d=1-b$  at each generation). Once a deme exceeds the maximum size (10,000 cells), it splits into two offspring demes via random sampling of cells from a binomial distribution ( $N_c, 0.5$ ), where  $N_c$  is the current deme size.

To model monoclonal seeding, a single cell at the tumor periphery was randomly sampled as the metastasis founder cell. To model polyclonal seeding, a cluster of cells ( $n=10$ ) were randomly sampled from the whole tumor in order to maximize the clonal diversity within the metastasis founder cells. This is because if the clonal diversity in the metastasis founder cells is low, it essentially models the scenario of monoclonal seeding by a cluster of genetically similar cells. The metastasis grows at same spatial model with primary tumor started from the metastasis founder cell or cell cluster ( $n=10$ ). During each cell division in the growth of primary tumor and metastasis, the number of neutral passenger mutations acquired in the coding portion of the genome follows a Poisson distribution with mean  $u$ . Thus, the probability that  $k$  mutations occurred in each cell division is as follows:

$$P(x = k) = \frac{u^k e^{-u}}{k!} \quad \text{Eq. (4)}$$

where an infinite sites model and constant mutation rate are assumed during tumor progression. Advantageous mutations also arise stochastically via a Poisson process with mean  $u_s$  during each cell division. We assume  $u_s=10^{-5}$  per cell division in the genome and each increases the cell division probability<sup>70</sup>. The cell birth and death probabilities for a selectively beneficial clone are  $b_s=b \times (1+s)$  and  $d_s=1-d_s=1-b \times (1+s)$ , respectively, thus the selective advantage for an advantageous mutation is defined as  $s=b_s/b-1$ .

During simulation of primary and metastatic growth, each mutation is assigned a unique index that is recorded with respect to its genealogy and host cells, enabling analysis of the mutational frequency in a bulk sample of tumor cells during different stages of growth. We simulate growth until the primary and metastasis reach a size of  $\sim 10^9$  cells (or  $\sim 10 \text{ cm}^3$ ) and then sample a bulk subpopulation (consisting of  $\sim 10^6$  cells) at the peripheral region of the primary tumor and metastasis, respectively. The VAF of all SSNVs in the sampled bulk subpopulation is considered the true VAF (denoted by  $f_T$ ), whereas the observed allele frequency is obtained via a statistical model that mimics the random sampling of alleles during sequencing. Specifically, we employ a Binomial distribution ( $n, f_T$ ) to generate the observed VAF at each site given its true frequency  $f_T$  and number of covered reads  $n$ . The number of covered reads at each site is assumed to follow a negative-binomial distribution (*Negative Binomial(size, depth)*) where depth is the mean sequencing depth and size corresponds to the variation parameter. We assume *depth*=100 and *size*=2 for the sequencing data in each tumor region and tissue purity=0.6 in order to model normal cell contamination in clinical samples. A mutation is called when the number of variant reads is  $\geq 3$ , thereby applying the same criteria as for the patient tumors.

We employed a mutation rate  $u=0.6$  per cell division in the exonic region (corresponding to  $10^{-8}$  per site per cell division in the 60Mb diploid coding regions). In order to model varying scenarios of tumor growth dynamics, selection and timing of metastatic dissemination, for each primary tumor/metastasis (P/M) pair, the birth probability  $b$  of founding cells, selection coefficient  $s$  and primary tumor size at dissemination  $N_d$  was sampled from a uniform distribution,  $b \sim U(0.55, 0.65)$ ,  $\log_{10}(s) \sim U(-3, -1)$  and  $\log_{10}(N_d) \sim U(4, 8)$ , respectively. 500 virtual P/M pairs were simulated under each of the monoclonal seeding and polyclonal seeding scenarios. The number of M-private clonal SSNVs ( $L_m$ ), P-private clonal SSNVs ( $L_p$ ) and P/M shared subclonal SSNVs ( $W_s$ ) for each P/M pair were counted from the simulation data and the simulated JSI was computed by Eq.(3).

## Data availability

The exome sequencing data for in-house collected colorectal cancer patients have been deposited at the European Genotype Phenotype Archive (EGA) under accession number EGAS0000100XXXX. The accession numbers for public datasets were listed in Table S1.

## Code availability

Code used for genomic data analysis are available from: <https://github.com/cancersysbio>

## References

1. Talmadge, J.E., Wolman, S.R. & Fidler, I.J. Evidence for the clonal origin of spontaneous metastases. *Science* **217**, 361-3 (1982).
2. Yamamoto, N. *et al.* Determination of clonality of metastasis by cell-specific color-coded fluorescent-protein imaging. *Cancer Res* **63**, 7785-90 (2003).
3. Liu, W. *et al.* Copy number analysis indicates monoclonal origin of lethal metastatic prostate cancer. *Nat Med* **15**, 559-65 (2009).
4. Huang, Y. *et al.* Multilayered molecular profiling supported the monoclonal origin of metastatic renal cell carcinoma. *Int J Cancer* **135**, 78-87 (2014).
5. Gundem, G. *et al.* The evolutionary history of lethal metastatic prostate cancer. *Nature* **520**, 353-357 (2015).
6. Maddipati, R. & Stanger, B.Z. Pancreatic Cancer Metastases Harbor Evidence of Polyclonality. *Cancer Discov* **5**, 1086-97 (2015).
7. Cheung, K.J. *et al.* Polyclonal breast cancer metastases arise from collective dissemination of keratin 14-expressing tumor cell clusters. *Proc Natl Acad Sci U S A* **113**, E854-63 (2016).
8. Hunter, K.W., Amin, R., Deasy, S., Ha, N.H. & Wakefield, L. Genetic insights into the morass of metastatic heterogeneity. *Nat Rev Cancer* **18**, 211-223 (2018).
9. Klein, C.A. Parallel progression of primary tumours and metastases. *Nat Rev Cancer* **9**, 302-12 (2009).
10. Naxerova, K. & Jain, R.K. Using tumour phylogenetics to identify the roots of metastasis in humans. *Nat Rev Clin Oncol* **12**, 258-72 (2015).
11. Robinson, D.R. *et al.* Integrative clinical genomics of metastatic cancer. *Nature* **548**, 297-303 (2017).
12. Bertucci, F. *et al.* Genomic characterization of metastatic breast cancers. *Nature* **569**, 560-564 (2019).
13. Priestley, P. *et al.* Pan-cancer whole-genome analyses of metastatic solid tumours. *Nature* (2019).
14. Zhao, Z.M. *et al.* Early and multiple origins of metastatic lineages within primary tumors. *Proc Natl Acad Sci U S A* **113**, 2140-5 (2016).
15. Macintyre, G. *et al.* How Subclonal Modeling Is Changing the Metastatic Paradigm. *Clin Cancer Res* **23**, 630-635 (2017).
16. Hu, Z. *et al.* Quantitative evidence for early metastatic seeding in colorectal cancer. *Nat Genet* **51**, 1113-1122 (2019).
17. Leung, M.L. *et al.* Single-cell DNA sequencing reveals a late-dissemination model in metastatic colorectal cancer. *Genome Res* **27**, 1287-1299 (2017).

18. Turajlic, S. & Swanton, C. Metastasis as an evolutionary process. *Science* **352**, 169-75 (2016).
19. Lee, S.Y. *et al.* Comparative genomic analysis of primary and synchronous metastatic colorectal cancers. *PLoS One* **9**, e90459 (2014).
20. Kim, T.M. *et al.* Subclonal Genomic Architectures of Primary and Metastatic Colorectal Cancer Based on Intratumoral Genetic Heterogeneity. *Clin Cancer Res* **21**, 4461-72 (2015).
21. Lim, B. *et al.* Genome-wide mutation profiles of colorectal tumors and associated liver metastases at the exome and transcriptome levels. *Oncotarget* **6**, 22179-90 (2015).
22. Uchi, R. *et al.* Integrated Multiregional Analysis Proposing a New Model of Colorectal Cancer Evolution. *PLoS Genet* **12**, e1005778 (2016).
23. Brastianos, P.K. *et al.* Genomic Characterization of Brain Metastases Reveals Branched Evolution and Potential Therapeutic Targets. *Cancer Discov* **5**, 1164-1177 (2015).
24. Um, S.W. *et al.* Molecular Evolution Patterns in Metastatic Lymph Nodes Reflect the Differential Treatment Response of Advanced Primary Lung Cancer. *Cancer Res* **76**, 6568-6576 (2016).
25. Chung, W. *et al.* Single-cell RNA-seq enables comprehensive tumour and immune cell profiling in primary breast cancer. *Nat Commun* **8**, 15081 (2017).
26. Ng, C.K.Y. *et al.* Genetic Heterogeneity in Therapy-Naive Synchronous Primary Breast Cancers and Their Metastases. *Clin Cancer Res* **23**, 4402-4415 (2017).
27. Razavi, P. *et al.* The Genomic Landscape of Endocrine-Resistant Advanced Breast Cancers. *Cancer Cell* **34**, 427-438 e6 (2018).
28. Siegel, M.B. *et al.* Integrated RNA and DNA sequencing reveals early drivers of metastatic breast cancer. *J Clin Invest* **128**, 1371-1383 (2018).
29. Ullah, I. *et al.* Evolutionary history of metastatic breast cancer reveals minimal seeding from axillary lymph nodes. *J Clin Invest* **128**, 1355-1370 (2018).
30. Sun, R. *et al.* Between-region genetic divergence reflects the mode and tempo of tumor evolution. *Nat Genet* **49**, 1015-1024 (2017).
31. Curtis, C. *et al.* The genomic and transcriptomic architecture of 2,000 breast tumours reveals novel subgroups. *Nature* **486**, 346-52 (2012).
32. Pich, O. *et al.* The mutational footprints of cancer therapies. *bioRxiv*, 683268 (2019).
33. Christensen, S. *et al.* 5-Fluorouracil treatment induces characteristic T>G mutations in human cancer. *Nat Commun* **10**, 4571 (2019).
34. Adzhubei, I.A. *et al.* A method and server for predicting damaging missense mutations. *Nat Methods* **7**, 248-9 (2010).
35. Rogers, M.F. *et al.* FATHMM-XF: accurate prediction of pathogenic point mutations via extended features. *Bioinformatics* **34**, 511-513 (2018).
36. Tokheim, C. & Karchin, R. CHASMplus Reveals the Scope of Somatic Missense Mutations Driving Human Cancers. *Cell Syst* **9**, 9-23 e8 (2019).
37. Patel, S.A. & Vanharanta, S. Epigenetic determinants of metastasis. *Mol Oncol* **11**, 79-96 (2017).
38. Martincorena, I. *et al.* Universal Patterns of Selection in Cancer and Somatic Tissues. *Cell* **171**, 1029-1041 e21 (2017).

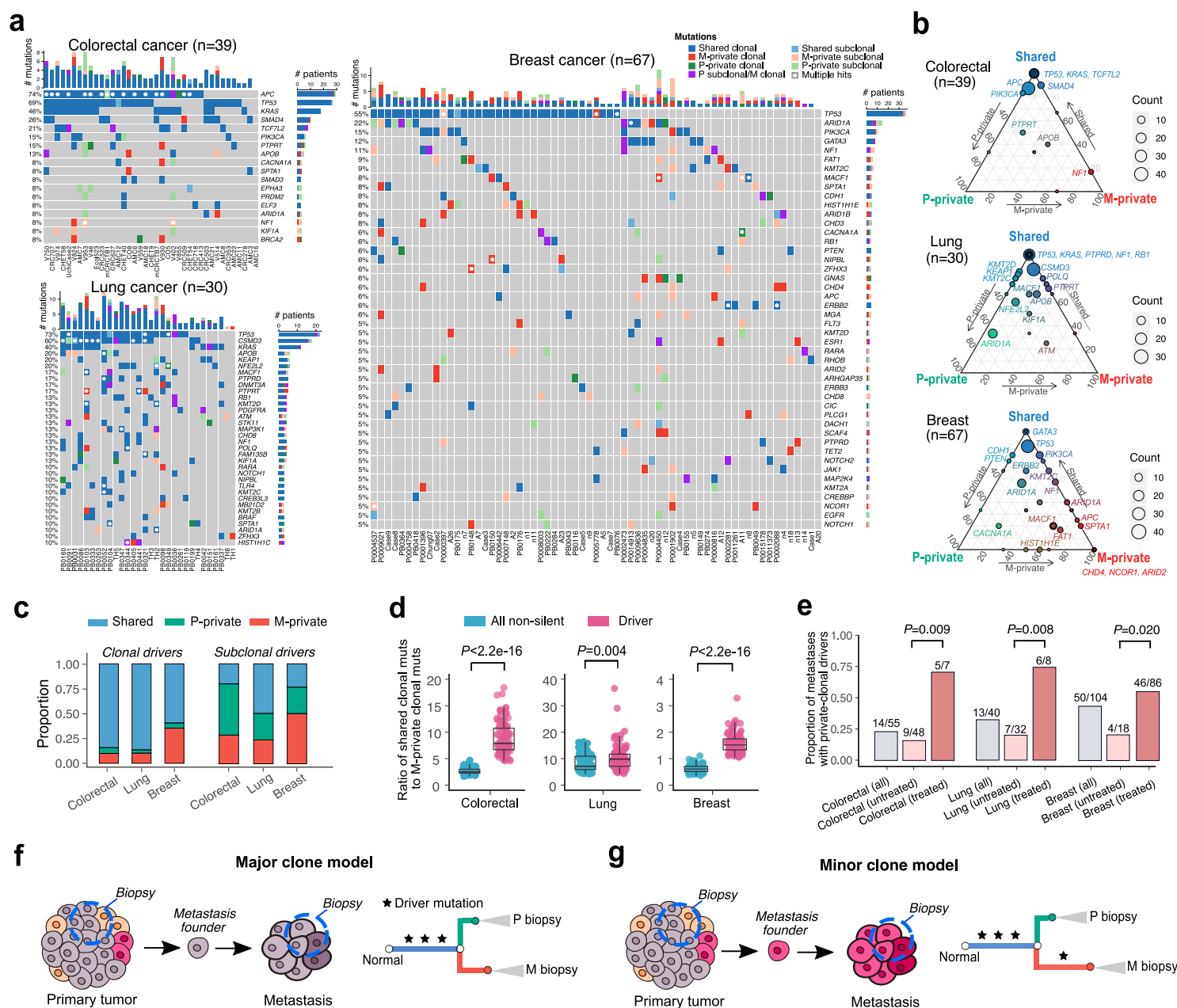
39. Makohon-Moore, A.P. *et al.* Limited heterogeneity of known driver gene mutations among the metastases of individual patients with pancreatic cancer. *Nat Genet* **49**, 358-366 (2017).
40. Sottoriva, A. *et al.* A Big Bang model of human colorectal tumor growth. *Nat Genet* **47**, 209-16 (2015).
41. Kang, H. *et al.* Many private mutations originate from the first few divisions of a human colorectal adenoma. *J Pathol* **237**, 355-62 (2015).
42. Williams, M.J. *et al.* Quantification of subclonal selection in cancer from bulk sequencing data. *Nat Genet* **50**, 895-903 (2018).
43. Williams, M.J., Werner, B., Barnes, C.P., Graham, T.A. & Sottoriva, A. Identification of neutral tumor evolution across cancer types. *Nat Genet* **48**, 238-244 (2016).
44. Benzekry, S. *et al.* Classical mathematical models for description and prediction of experimental tumor growth. *PLoS Comput Biol* **10**, e1003800 (2014).
45. Stein, R.G. *et al.* The impact of breast cancer biological subtyping on tumor size assessment by ultrasound and mammography - a retrospective multicenter cohort study of 6543 primary breast cancer patients. *BMC Cancer* **16**, 459 (2016).
46. Rueda, O.M. *et al.* Dynamics of breast-cancer relapse reveal late-recurring ER-positive genomic subgroups. *Nature* **567**, 399-404 (2019).
47. Harper, K.L. *et al.* Mechanism of early dissemination and metastasis in Her2(+) mammary cancer. *Nature* (2016).
48. Hosseini, H. *et al.* Early dissemination seeds metastasis in breast cancer. *Nature* (2016).
49. Reiter, J.G. *et al.* Minimal functional driver gene heterogeneity among untreated metastases. *Science* **361**, 1033-1037 (2018).
50. Aguirre-Ghiso, J.A. Models, mechanisms and clinical evidence for cancer dormancy. *Nat Rev Cancer* **7**, 834-46 (2007).
51. Paez, D. *et al.* Cancer dormancy: a model of early dissemination and late cancer recurrence. *Clin Cancer Res* **18**, 645-53 (2012).
52. Recasens, A. & Munoz, L. Targeting Cancer Cell Dormancy. *Trends Pharmacol Sci* **40**, 128-141 (2019).
53. Zijlstra, A. *et al.* The importance of developing therapies targeting the biological spectrum of metastatic disease. *Clin Exp Metastasis* (2019).

## References in Methods

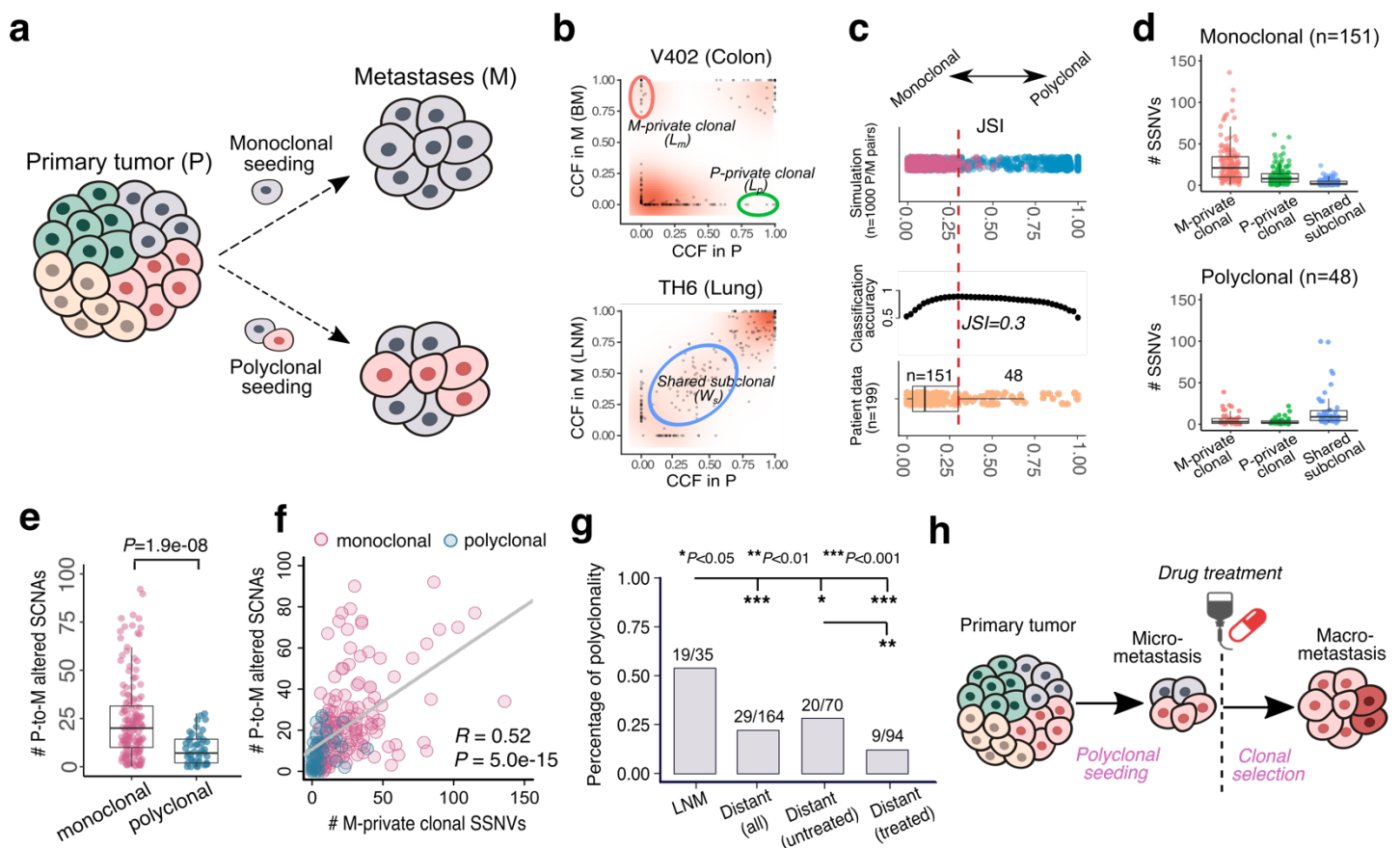
54. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754-60 (2009).
55. Cibulskis, K. *et al.* Sensitive detection of somatic point mutations in impure and heterogeneous cancer samples. *Nat Biotechnol* **31**, 213-9 (2013).
56. Koboldt, D.C. *et al.* VarScan 2: somatic mutation and copy number alteration discovery in cancer by exome sequencing. *Genome Res* **22**, 568-76 (2012).
57. Costello, M. *et al.* Discovery and characterization of artifactual mutations in deep coverage targeted capture sequencing data due to oxidative DNA damage during sample preparation. *Nucleic Acids Res* **41**, e67 (2013).
58. Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078-9 (2009).



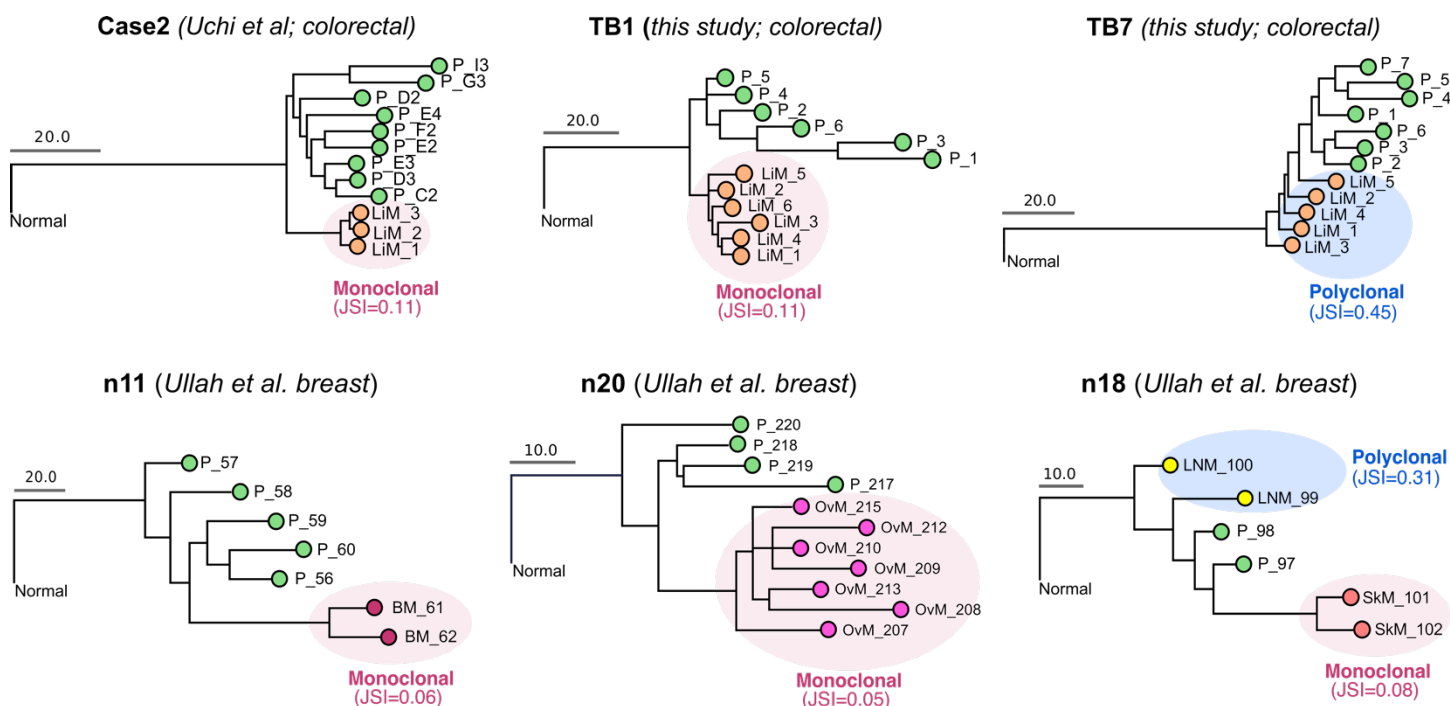
59. Saunders, C.T. *et al.* Strelka: accurate somatic small-variant calling from sequenced tumor-normal sample pairs. *Bioinformatics* **28**, 1811-7 (2012).
60. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic variants from high-throughput sequencing data. *Nucleic Acids Res* **38**, e164 (2010).
61. Ha, G. *et al.* TITAN: inference of copy number architectures in clonal cell populations from tumor whole-genome sequence data. *Genome Res* **24**, 1881-93 (2014).
62. Ha, G. *et al.* Integrative analysis of genome-wide loss of heterozygosity and monoallelic expression at nucleotide resolution reveals disrupted pathways in triple-negative breast cancer. *Genome Res* **22**, 1995-2007 (2012).
63. Li, B. & Li, J.Z. A general framework for analyzing tumor subclonality using SNP array and DNA sequencing data. *Genome Biol* **15**, 473 (2014).
64. McGranahan, N. *et al.* Clonal status of actionable driver events and the timing of mutational processes in cancer evolution. *Sci Transl Med* **7**, 283ra54 (2015).
65. Bailey, M.H. *et al.* Comprehensive Characterization of Cancer Driver Genes and Mutations. *Cell* **174**, 1034-1035 (2018).
66. Zhou, Y. *et al.* Metascape provides a biologist-oriented resource for the analysis of systems-level datasets. *Nat Commun* **10**, 1523 (2019).
67. Diaz-Gay, M. *et al.* Mutational Signatures in Cancer (MuSiCa): a web application to implement mutational signatures analysis in cancer samples. *BMC Bioinformatics* **19**, 224 (2018).
68. Alexandrov, L.B. *et al.* Signatures of mutational processes in human cancer. *Nature* **500**, 415-21 (2013).
69. J., F. PHYLIP-phylogeny inference package (version 3.2). *cladistics* **5**, 6 (1989).
70. Bozic, I. *et al.* Accumulation of driver and passenger mutations during tumor progression. *Proc Natl Acad Sci U S A* **107**, 18545-50 (2010).



**Figure 1. Landscape of driver mutations in paired primary tumors (P) and metastases (M).** (a) Oncoprint of functional driver mutations in the three cancer types grouped by P/M shared, P-private or M-private mutations including both clonal or subclonal drivers. Genes mutated in at least three patients are shown. Boxes with white circles indicate genes with multiple mutations in a given patient usually in tumor suppressor genes (*TP53*, *APC*, *CSMD3*, etc.). (b) Ternary plot of mutation counts in driver genes, comparing P-private (left, green), M-private (right, red), and shared (top, blue). The color of each circle indicates the relative frequency of driver mutations among these groups, while the size of the circle represents their overall count in the corresponding cancer type. (c) The proportion of shared, P-private or M-private drivers (*clonal*: shared clonal, P-private clonal or M-private clonal; *subclonal*: shared subclonal, P subclonal/M clonal, P-private subclonal or M-private subclonal) in each of the three cancer types. (d) The ratio of shared clonal to M-private clonal mutations for all non-silent and driver mutations, respectively. A down-sampling procedure was performed to derive the ratio (Methods) where n=100 down-samplings (50% patients each) were repeated for each of the three cancer types. *P*-value, Wilcoxon Rank-Sum Test (two-sided). Bar, median; box, 25th to 75th percentile (interquartile range, IQR); vertical line, data within 1.5 times the IQR. (e) The proportion of metastases harboring at least one private clonal driver mutation grouped by all metastases, untreated and treated metastases. *P*-value, Fisher's exact test (two-sided). (f) Schematic representation of the major clone model where metastasis originates from the major driver clone in the primary tumor leading to driver gene homogeneity between paired P and M biopsies. (g) Schematic representation of the minor clone model in which metastases originate from a minor clone in the primary tumor. Due to the inability to detect the minor driver clone in bulk sequencing data, the minor clone model leads to driver heterogeneity between P and M biopsies.

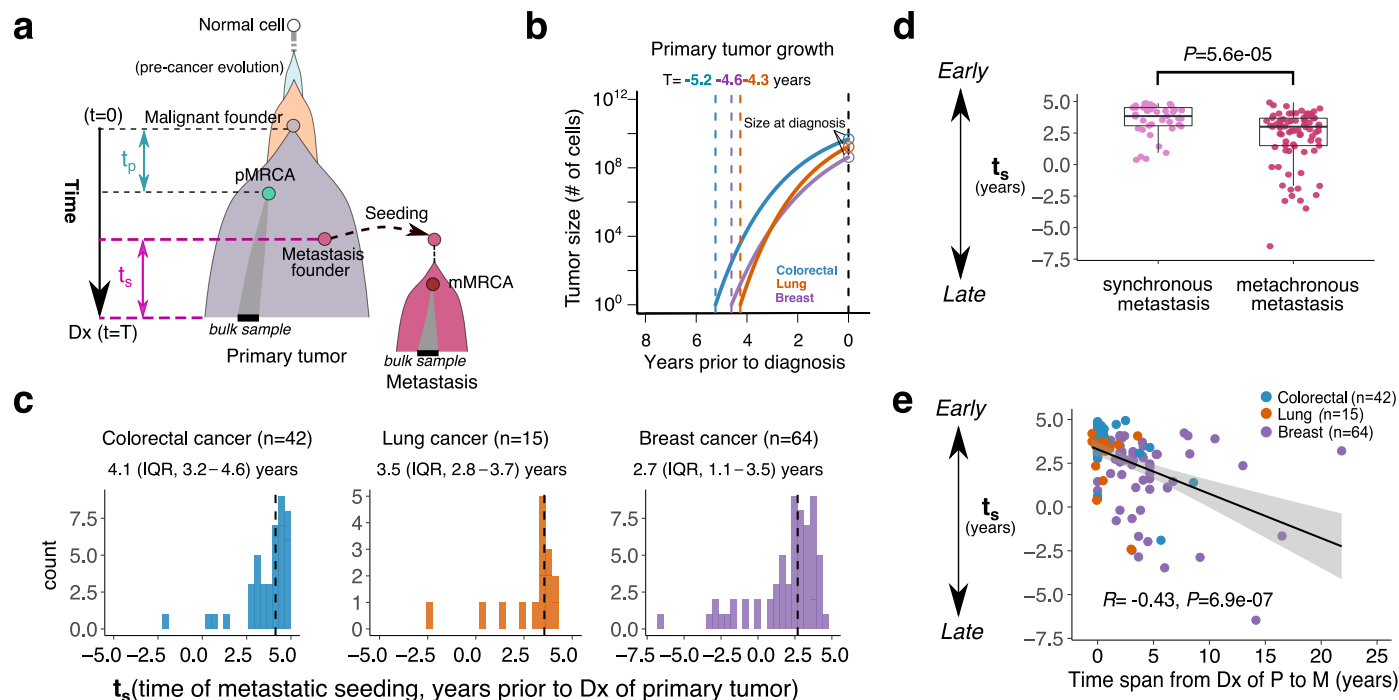


**Figure 2. The clonality of lymph node and distant metastases.** (a) Schematic illustration of monoclonal versus polyclonal seeding for a single metastasis. (b) Distinct patterns of monoclonal versus polyclonal seeding based on the cancer cell fraction (CCF) of SSNVs between P/M pairs. An example patient is shown for each scenario: monoclonal seeding (colon cancer patient V402 with brain metastasis (BM)); polyclonal seeding (lung cancer patient TH6 with lymph node metastasis (LNM)). Green and red circles indicate the P-private clonal SSNVs (the number denoted by  $L_p$ ) and M-private clonal SSNVs (the number denoted by  $L_m$ ), respectively. Blue circle indicates the P/M shared subclonal SSNVs (the number denoted by  $W_s$ ). (c) Classification of monoclonal versus polyclonal seeding based on the Jaccard similarity index (JSI). Top, JSI values in 1000 virtual P/M tumor pairs simulated from a spatial tumor growth model in which 500 were from monoclonal seeding (number of metastasis founder cell=1) and 500 were from polyclonal seeding (number of metastasis founder cells=10). Middle, classification accuracy by varying the cutoff of JSI from 0 to 1 based on the simulation data. Bottom, the JSI values in patient data (n=199 P/M pairs) where the cutoff JSI=0.3 was used to identify monoclonal seeding (n=151) or polyclonal seeding (n=48). (d)  $L_m$ ,  $L_p$ ,  $W_s$  values in the patient data. Top, monoclonal metastases; bottom, polyclonal metastases. Bar, median; box, 25th to 75th percentile (interquartile range, IQR); vertical line, data within 1.5 times the IQR. (e) The number of P-to-M altered SCNAs for monoclonal and polyclonal metastases, respectively. (f) Positive correlation between  $L_m$  and the number of P-to-M altered SCNAs. n=199 P/M pairs and Pearson's correlation ( $R$ ) and  $P$ -value were reported. (g) Polyclonal seeding is common in lymph node metastases (LNM) and untreated distant metastases relative to treated distant metastases. (h) Schematic illustration of the scenario where treatment promotes monoclonality as a result of selection for a resistant subclone, despite initial seeded by polyclonal disseminated cells.

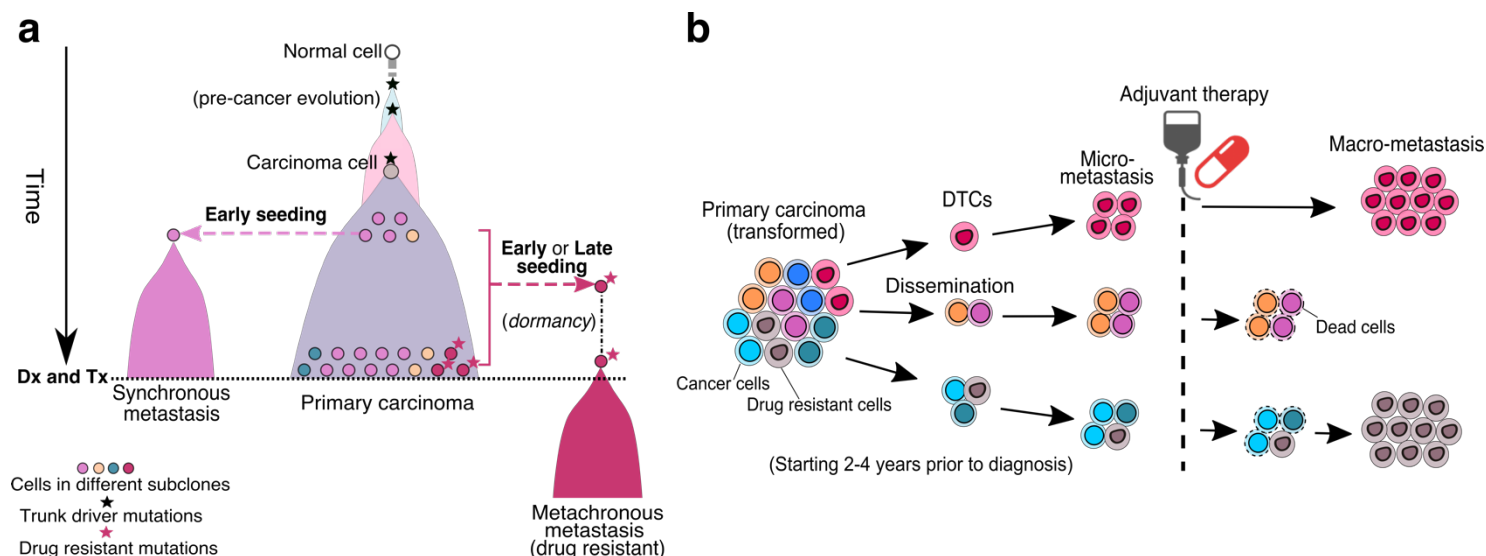


**Figure 3. Tumor sample phylogenies based on multi-region sequencing data.** The maximum parsimony method was used to reconstruct multi-sample trees for each patient based on the presence or absence SSNVs/indels amongst the samples while accounting for the loss-of-heterozygosity in the mutant sites. For each P/M sample pair, the Jaccard similarity index (JSI) was computed according to Eq. (3) based on the numbers of M-private clonal, P-private clonal and P-M shared subclonal SSNVs. High JSI values ( $>0.3$ ) indicates polyclonal seeding while low JSI values ( $\leq 0.3$ ) indicates monoclonal seeding. Monoclonal seeding gives rise to monophyletic tree structures (pink shading indicates metastatic samples within a single phylogenetic clade), whereas polyclonal seeding gives rise to a polyphyletic structure (blue shading indicates metastatic samples within multiple phylogenetic clades) in the metastasis samples. P, primary tumor; OvM, ovarian metastasis; LNM, lymph node metastasis; SkM, skin metastasis; LiM, liver metastasis. Additional patient data are shown in **Fig. S13**.





**Figure 4. Chronology of metastatic seeding.** (a) Schematic for the timing of metastatic seeding prior to diagnosis of the primary tumor in number of years,  $t_s$ .  $T$  denotes the total time of primary tumor expansion from emergence of the malignant founder cell to diagnosis while  $t_p$  denotes the time from emergence of the malignant founder cell to the most recent common ancestor (MRCA) of cells in primary bulk sample (denoted pMRCA).  $t_s$  can be estimated by Eq.(1). Dx, diagnosis (b) Estimation of the average  $T$  with a Gompertzian growth model is 5.2 (interquartile range or IQR, 4.3–7.7), 4.3 (IQR, 2.7–4.4) and 4.6 (IQR, 3.2–6.6) years for colorectal, lung and breast cancer, respectively. (c) Estimation of the time of metastatic seeding ( $t_s$ ) for individual distant metastases (monoclonal) in each cancer types. The median  $t_s$  and IQR are shown. Negative  $t_s$  indicates that the metastasis was seeded after the diagnosis of primary tumor. (d) The distribution of  $t_s$  in synchronous metastases ( $n=40$ ) and metachronous metastases ( $n=81$ ).  $P$ -value, Wilcoxon Rank-Sum Test (two-sided). Bar, median; box, 25th to 75th percentile (IQR); vertical line, data within 1.5 times the IQR. (e) Correlation between  $t_s$  and the time span from diagnosis of primary tumor to metastasis. Pearson's correlation ( $R$ ) and  $P$ -value are reported.



**Figure 5. Schematic model of metastatic spread and the impact of therapy.** (a) Schematic illustration of early versus late metastatic seeding leading to synchronous and metachronous metastases. Metastatic seeding starts quickly following the emergence of founding carcinoma cell. Synchronous metastasis, which exhibits low genomic divergence with primary tumor, is seeded early by the major clone in primary tumor. Metachronous metastasis, exhibit higher genomic divergence relative to the primary tumor and often emerge after adjuvant therapy. Metachronous metastasis can be seeded either early or late depending on selective pressure by treatment and/or latency period of dormant disseminated cells<sup>50,51</sup>. Metachronous metastases with specific driver mutations that confer resistance can be selected leading to high genomic divergence between the primary tumor and treated metastasis. Dx, diagnosis; Tx, treatment. (b) Treatment (here adjuvant therapy) remodels the clonal architecture of metastasis. Dissemination and metastatic seeding (monoclonal or polyclonal) initially gives rise to undetectable micro-metastases. While treatment may eliminate drug-sensitive lesions, those that are resistant grow out. Metastatic relapse following adjuvant treatment may be delayed by treatment, but this may result in a more aggressive, resistant lesion. DTCs, disseminated tumor cells.