

# 1 The human pathome shows sex specific 2 aging patterns post- development

3 Michael Ben Ezra<sup>1, 2, \*</sup> Jonas Bach Garbrecht<sup>1</sup>, Nasya Rasmussen<sup>1</sup>, Indra  
4 Heckenbach<sup>1, 3</sup>, Michael A. Petr<sup>1, 3</sup>, Daniela Bakula<sup>1</sup>, Laust Mortensen<sup>2</sup> & Morten  
5 Scheibye-Knudsen<sup>1, 3, \*</sup>

6  
7 <sup>1</sup>Center for Healthy Aging, Department of Cellular and Molecular Medicine,  
8 University of Copenhagen, <sup>2</sup>Methods and Analysis, Statistics Denmark, <sup>3</sup>Tracked.bio  
9 \*Correspondence: [mbenezra@sund.ku.dk](mailto:mbenezra@sund.ku.dk), [mscheibye@sund.ku.dk](mailto:mscheibye@sund.ku.dk)

## 10 Abstract

11 Little is known about tissue specific changes that occur with aging in humans. Using  
12 the description of 33 million histological samples we extract thousands of age- and  
13 mortality-associated features from text narratives that we call The Human Pathome  
14 (pathoage.com). Notably, we can broadly determine when pathological aging starts,  
15 indicating a sexual dimorphism with females aging earlier but slower and males  
16 aging later but faster. Using machine learning, we employ unsupervised topic-  
17 modelling to identify terms and themes that predict age and mortality. As a proof of  
18 principle, we cross reference these terms in PubMed to identify nintedanib as a  
19 potential aging intervention and show that nintedanib reduces markers of cellular  
20 senescence, reduces pro-fibrotic gene pathways in senescent cells and extends the  
21 lifespan of fruit flies. Our findings pave the way for expanded exploitation of  
22 population datasets towards discovery of novel aging interventions.

## 24 Introduction

25 Aging is a complex, multifactorial process<sup>1,2</sup> that leads to declining physiology and a  
26 susceptibility to disease<sup>3</sup>. Yet, little is known about tissue specific changes that occur  
27 with aging in humans. Clinical text constitutes the most abundant data type in  
28 electronic health care records which are implemented in most countries<sup>4</sup>.  
29 Specifically, pathology records are rich in descriptions of cellular and histological  
30 samples of healthy and diseased human tissue and therefore represent a  
31 considerable opportunity to systematically characterize tissue specific changes that  
32 occur in aging. Nonetheless, electronic health care records are a vastly underused  
33 data resource due to their limited availability to researchers<sup>5</sup>. Furthermore,  
34 unstructured text data are not directly amenable to computational analysis and  
35 clinical text is highly heterogeneous. Importantly, using natural language processing  
36 and machine learning, phenotypes can be extracted from clinical text and used to  
37 discover correlations and stratify patient cohorts<sup>6,7</sup>.

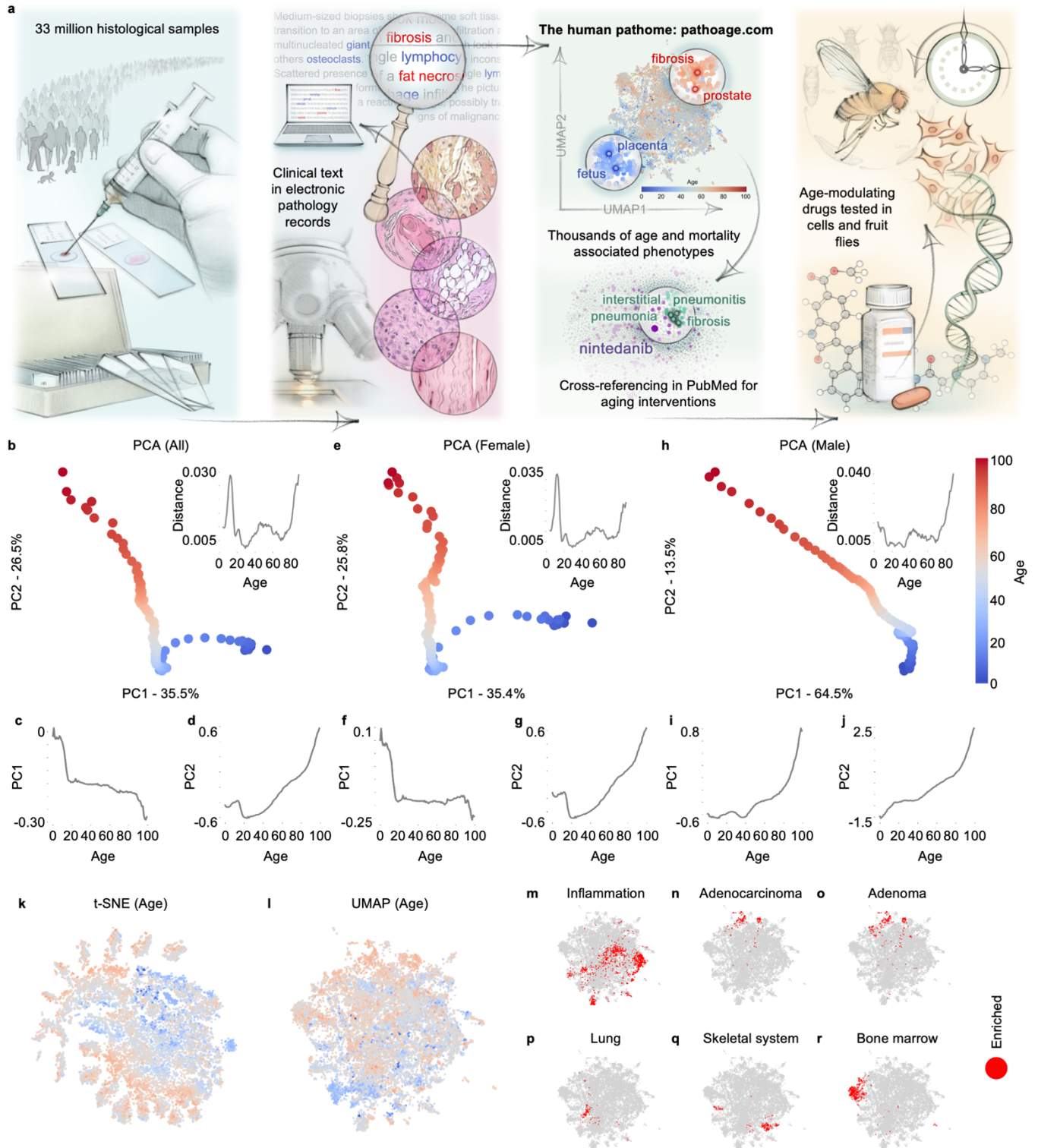
38  
39 To get an unbiased description of organismal- and tissue-specific aging, we  
40 analyzed the Danish pathology register containing the clinical description of over 33  
41 million samples collected since 1970<sup>8</sup> from over 4.9 million individuals some born as  
42 early as 1876 (Fig. 1a). Using natural language processing we extracted thousands  
43 of clinical features from unstructured pathology narrative texts. We combine this with  
44 vital statistics to identify age- and mortality-associated features. Using supervised  
45 and unsupervised machine-learning we identify population-based patterns of aging  
46 and surprisingly discover that pathological aging starts almost immediately after  
47 development in the late teens for females. For males, pathological aging starts later

48 (~40 years) but progresses faster. Conversely, tissue-specific patterns of aging show  
49 that some tissues age linearly and others age along developmental and pathological  
50 aging trajectories. To further investigate the meaning of clinical features we  
51 employed topic-modelling<sup>9</sup> and reveal specific age- and mortality associated themes.  
52 As a proof of principle, we deploy this in lung pathology records and find that the  
53 predicative power of topic modelling themes is stronger than individual features. We  
54 further cross-reference the age-associated terms from the pathology datasets within  
55 all published PubMed abstracts and identify compounds enriched in aging terms.  
56 Among them, we identify nintedanib, a tyrosine kinase inhibitor<sup>10</sup>, as a potential  
57 pharmacological intervention in aging. Indeed, nintedanib, an antifibrotic agent,  
58 reduces markers of cellular senescence, reduces pro-fibrotic gene pathways in  
59 senescent cells and extend the lifespan of *drosophila melanogaster*.

## 60 Results

### 61 Pathological aging begins post-development for females and at 62 mid-life for males

63 To explain the variance in pathology records in the entire pathology register we  
64 identified the average term frequency within each age-group (0-100) and performed  
65 a principal component analysis. Remarkably, we observed a strong correlation  
66 between the main principal components PC1 (35.5%) and PC2 (26.52%) and age,  
67 showing that variance in pathology records is strongly explained by age (Fig. 1b).  
68 We observed that the ages of development (0-18) vary primarily along PC1 (Fig. 1c)  
69 whereas post development ages 19 and over vary primarily along PC2 (Fig. 1d)



**Fig. 1 | Pathological aging begins post-development for females and at mid-life for males.** **a**, The Human Aging Pathome and aging intervention discovery concept and workflow. **b**, PCA of age-aggregated pathology records ( $n=20,316,270$ ) from in entire pathology register. Normalized Euclidean distance between age adjacent PCA coordinates. **c, d**, PC1, PC2 coordinates vs. Age. **e**, PCA of age-aggregated of pathology records ( $n=14,492,989$ ) from females in the entire pathology register. Normalized Euclidean distance between age adjacent PCA coordinates. **f, g**, PC1, PC2 coordinates vs. Age. **h**, PCA of age-aggregated of pathology records ( $n=5,823,281$ ) from males in the entire pathology register. Normalized Euclidean distance between age adjacent PCA coordinates. **i, j**, PC1, PC2 coordinates vs. Age. **k**, t-SNE of clinical features in age-aggregated pathology records in the entire pathology register. **l**, UMAP of clinical features in pathology records in the entire pathology register. **m-o**, Positive morphology specific enrichment: Inflammation, Adenocarcinoma and Adenoma. **p-r**, Positive enrichment of clinical terms in tissue specific pathology records: Lung, Skeletal system and Bone marrow.

71 suggesting that PC1 describes variance in development while PC2 describes true  
72 pathological aging-associated changes. We also noted the Euclidean distance of age  
73 adjacent PCA components to assess the increase in variance with age. We noted a  
74 peak around the end of development, at midlife and late in life. Since we saw  
75 increased variance around midlife, we speculated that there could be sex-dependent  
76 differences in pathological aging perhaps around menopause. Strikingly, in females  
77 age-associated changes appear immediately after development (Fig. 1g) while in  
78 males pathological aging appears to start (Fig. 1i) at around forty years of age but  
79 does so at an increasing rate. Notably, while aging starts earlier in females the  
80 contribution of this factor to the overall variance is much smaller in females than  
81 males (PC2 females accounts for 25.8% of variance while PC1 for males account for  
82 64.5%). To visualize the entire landscape of terms, we applied t-distributed  
83 stochastic neighbor embedding (t-SNE) (Fig. 1k; Extended Data Fig. 1a,b for sex-  
84 specific t-SNEs) to the average term frequencies within each age-group overlaid with  
85 the mean incidence age of each feature in the pathology register. Strikingly, the t-  
86 SNE visualization shows that terms primarily associated with younger age groups  
87 coalesce in the center while terms associated with older-age groups project  
88 outwards in all directions reflecting an apparent age-dependent progression from  
89 order to disorder.

## 90 Patterns of tissue specific vocabulary identified in pathology 91 records

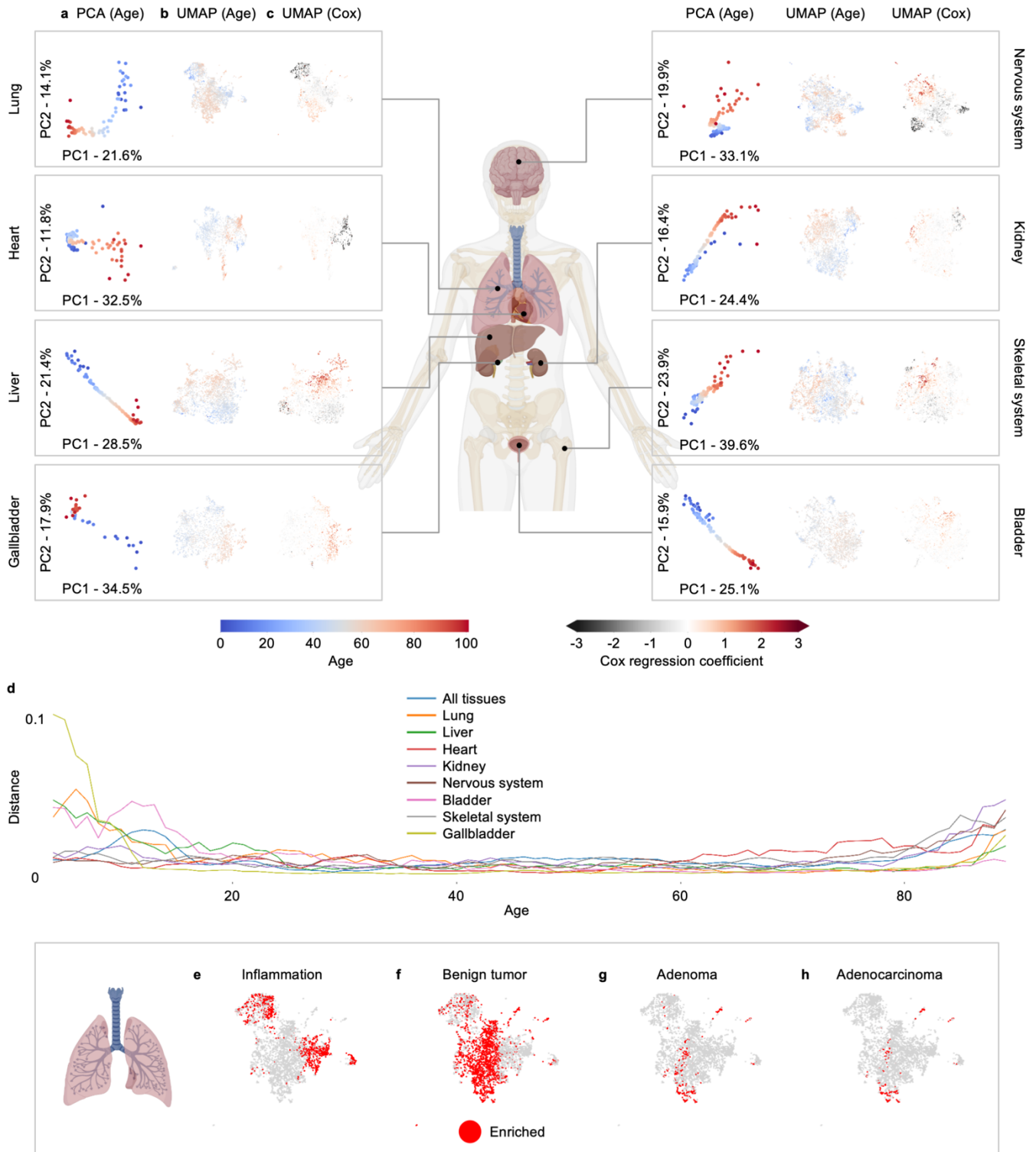
92 To visualize the co-occurrence of terms in the entire pathology register we applied  
93 Uniform Manifold Approximation and Projection (UMAP). The UMAP shows a more  
94 unidirectional age-effect than observed in the t-SNE (Fig. 1l; Extended Data Fig. 1c,d

95 for sex-specific UMAPs). In both t-SNE and UMAP visualizations we observed a  
96 tendency for terms with similar mean incidence age to co-occur. Pathology records  
97 in the registry are classified according to morphology and tissue (Fig. 1m-r; Extended  
98 Data Fig. 1e for additional tissues). We noted that records annotated with the  
99 morphology code 'inflammation' are enriched with terms from broad clusters of the  
100 feature landscape (Fig. 1m). Further, records associated with specific tissues such  
101 as lung, skeletal system and bone marrow are enriched in clinical terms from  
102 narrower regions of the feature landscape (Fig. 1p-r). Notably, terms enriched in lung  
103 tissue coincide with inflammation supporting the notion that inflammation affects lung  
104 more strongly compared with some other tissues in the body<sup>11</sup>. In addition, terms  
105 enriched in most tissues are largely non-overlapping, suggesting that terms used to  
106 describe specific tissues tend to be more distinct. On the other hand, related tissues  
107 such as skeletal system (Fig. 1q) and bone marrow (Fig. 1r) are enriched in terms  
108 from neighboring regions of the landscape. Altogether, these findings suggest that  
109 tissue specific patterns of aging can be identified in the dataset.

## 110 Tissues age along specific trajectories

111 To identify tissue specific aging patterns, we repeated the above analyses for every  
112 tissue in the body. We noted a similarly strong correlation between the two main  
113 principal components PC1 and PC2 and age in multiple tissues (Fig. 2a). To  
114 understand whether the mean incidence age of terms (Fig. 2b) and the mortality  
115 (defined as time to death from examination) associated with terms are correlated, we  
116 were able to connect the mortality data of over 1.3 million individuals with their own  
117 pathology records (Fig. 2c). Incidentally, hazard of death is associated with the word  
118 count length of clinical text narratives and with birth cohort (Extended Data Fig. 1f,g)





**Fig. 2 | Tissues age along specific trajectories.** **a**, PCA of age-aggregated tissue specific pathology records. **b**, UMAP of clinical features of tissue specific pathology records (mean incidence age). **c**, UMAP of clinical features of tissue specific pathology records (Cox regression coefficient). Tissues shown are: Lung (n=177,795), Liver (n=156,057), Heart (n=27,055), Kidney (n=85,244), Nervous system (n=183,729), Bladder (n=250,532), Skeletal system (n=242,282) and Gallbladder (n=182,261). **d**, Normalized Euclidean distance between age-adjacent PCA coordinates of all tissues in. **e-h**, Positive enrichment of clinical terms in morphology specific lung records: **e**, Inflammation. **f**, Benign tumor. **g**, Adenoma. **h**, Adenocarcinoma.

120 (Extended Data Fig. 2a-c for additional tissues). Indeed, age and mortality appear  
121 broadly correlated in all tissues. Interestingly, in several tissues (kidney, bladder,  
122 nervous system) we observe that age-related changes appear to be biphasic (Fig.  
123 2d) perhaps suggesting a phase associated with development and one associated  
124 with pathological aging. For other tissues, aging appears more linear (lung, liver,  
125 gallbladder, skeletal system, nervous system.) When investigating a single tissue  
126 such as lung, clinical features enriched in categories such as inflammation and  
127 benign tumors (Fig. 2e,f) appear to be mostly non-overlapping while adenoma and  
128 adenocarcinoma morphologies (Fig. 2g,h) appear to coincide. In sum, tissue specific  
129 trajectories define different patterns of aging indicating that different tissues age in  
130 different ways. To allow exploration of these phenomena, we have created a  
131 browsable database of the human pathome ([www.pathoage.com](http://www.pathoage.com)).

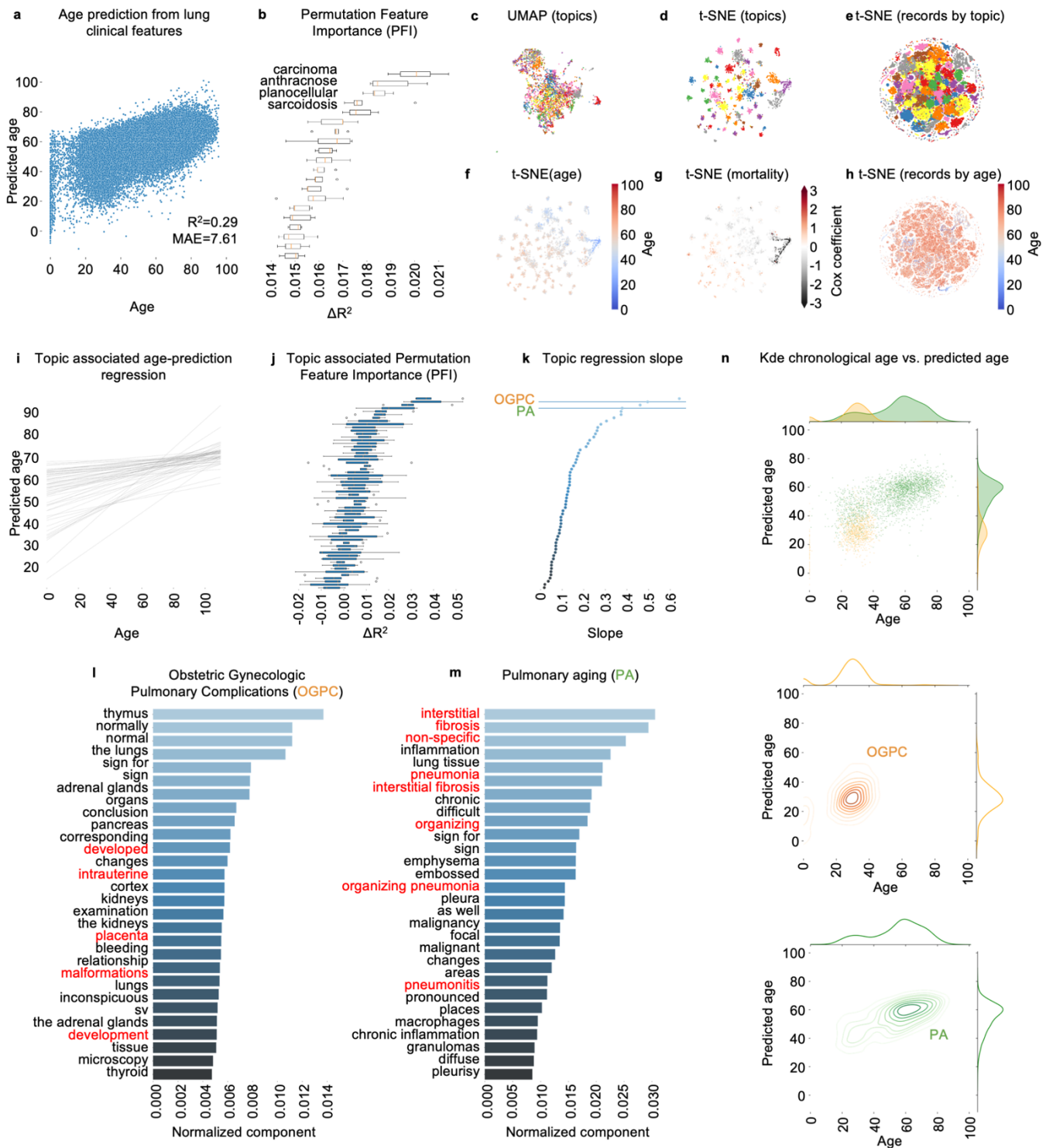
## 132 Clinical features in pathology text predict age

133 To better understand tissue specific aging, we fit a supervised deep neural network  
134 multilayer perceptron regression model to predict age from clinical text features in  
135 lung pathology records (Mean absolute error MAE=7.61, Fig. 3a). We performed a  
136 feature importance analysis (Fig. 3b) and identified the terms ‘carcinoma’,  
137 ‘anthracnose’, ‘planocellular’ and ‘sarcoidosis’ as most predictive of lung aging.  
138 However, given the relatively poor predictive power of the model ( $R^2=0.29$ ) and the  
139 small contribution of each term ( $\Delta R^2 < 0.021$ ) we decided to investigate whether a  
140 collection of associated terms would yield stronger predictive power.

141

142 To understand how clinical features semantically relate to one another we applied a  
143 latent Dirichlet allocation (LDA) topic model<sup>9</sup> to tissue-specific pathology records





**Fig. 3 | Semantic structures in clinical text describe lung pathologies and predict age.** **a**, DNN age-prediction from clinical features in lung pathology records. **b**, Permutation feature importance (PFI). **c**, UMAP of clinical features in lung records (LDA topic). **d**, t-SNE of clinical feature LDA distributions in topics (LDA topics). **e**, t-SNE of LDA record distributions in topics (LDA topics). **f**, t-SNE of clinical feature LDA distributions in topics (mortality: Cox regression coefficient). **g**, t-SNE of clinical feature LDA distributions in topics (age). **h**, t-SNE of LDA record distributions in topics (age). **i**, Topic associated age-prediction linear regression. **j**, Topic associated permutation feature importance (PFI). **k**, Topic associated age-prediction regression slope. **l**, Terms describing Obstetric Gynecologic Pulmonary Complications (OGPC). **m**, Terms describing pulmonary aging (PA). **n**, Bivariate kernel density estimation (kde) and histograms of chronological age vs. predicted age, records closely associated with PF and OGPC.

145 enabling us to identify clusters of co-occurring terms. We applied the topic model to  
146 177,795 lung pathology records and employed a model perplexity minimization  
147 strategy to determine that the clinical feature space is optimally decomposed into  
148 sixty topics (Extended Data Fig. 3a). A t-SNE visualization of clinical feature  
149 distributions in topics demonstrates the topic model's ability to segregate associated  
150 features into clusters (Fig. 3c,d). In turn, this also enables us to stratify individual  
151 pathology records by topic (Fig. 3e). Importantly, the topic model appears to identify  
152 collections of features with closely associated age and mortality (Fig. 3f,g; Extended  
153 Data Fig. 3b,c) suggesting that these semantic structures could describe clinically  
154 relevant themes. Importantly, stratified patient records also appear to have closely  
155 associated age at examination (Fig. 3h) further strengthening the notion that the  
156 topics we identified may characterize cohorts of similar individuals.

## 157 **Predicative power of topics is stronger than individual features**

158 To assess the predicative power of collections of associated terms, we performed  
159 linear regression on the age and predicted age of records closely associated with  
160 each topic (Fig. 3i; Extended Data Fig. 3d). Furthermore, we performed feature  
161 importance on the collected terms that make up each topic (Fig. 3j; Extended Data  
162 Fig. 3e). Notably, the maximum importance of collections of terms ( $\Delta R^2 < 0.039$ ) to  
163 age prediction is approximately 2-fold greater than that of the maximum importance  
164 of an individual feature ( $\Delta R^2 < 0.021$ ).

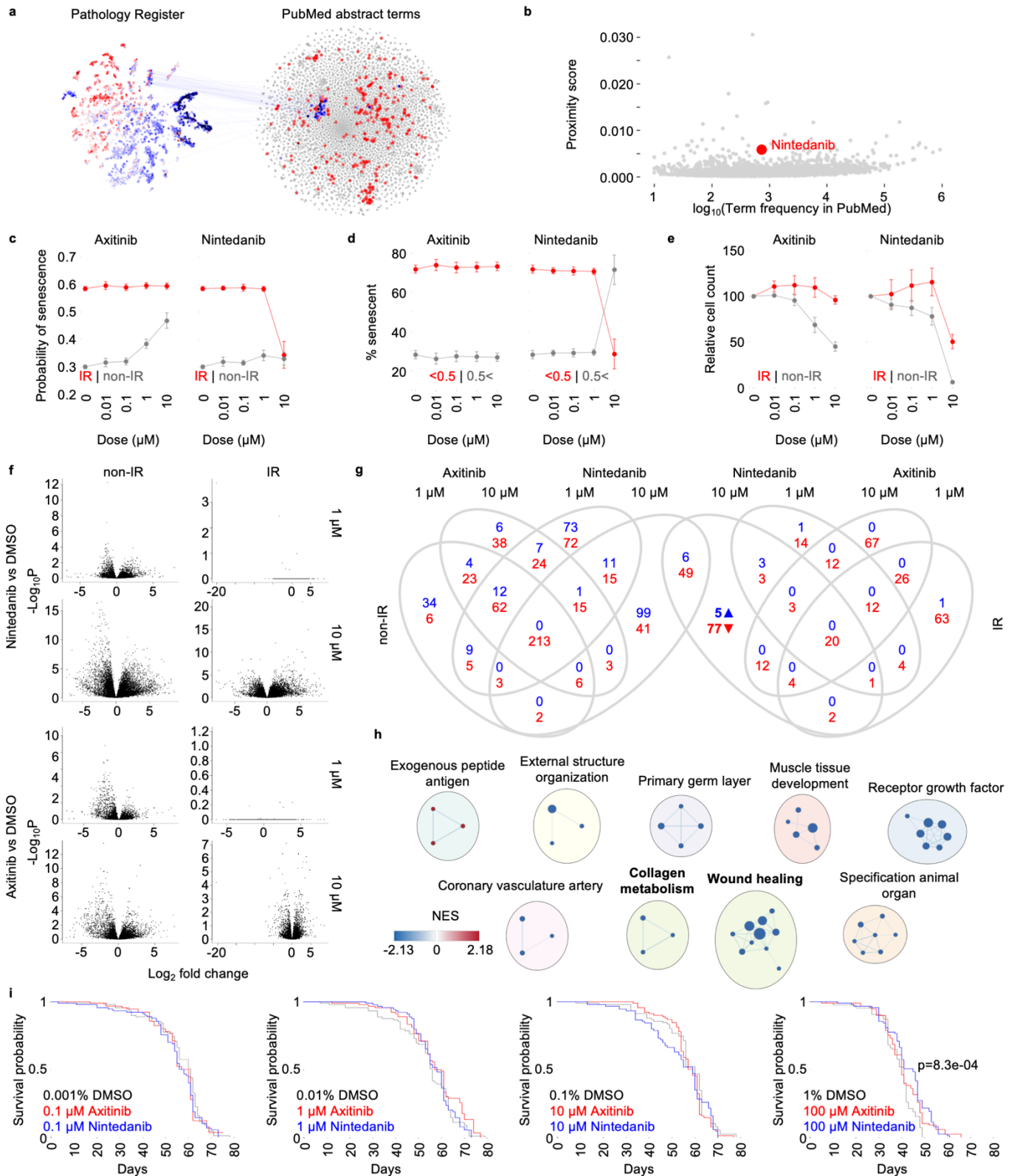
165

166 We then identified topics with changes in the age-prediction regression slope as  
167 topics where aging elicits alterations in the age-effect (Fig. 3k; Extended Data Fig.  
168 3f). Among the topics (Supplementary Table 2 for full list) we identified clinical

169 themes consisting of terms broadly describing cases of lung pathologies. One topic  
170 appeared to be associated with Human Immunodeficiency Virus (HIV) (terms such  
171 as 'fungi', 'pneumocystis', 'carinii', 'pneumocystis carinii', 'alveolar', 'inflammation'  
172 and 'fibrosis') (Extended Data Fig. 3g) and another with obstetric gynecologic  
173 pulmonary complications (OGPC) (terms such as 'development', 'intrauterine',  
174 'placenta' and 'malformations') (Fig. 3l). Interestingly, a topic appeared to describe  
175 pulmonary aging (PA) with terms such as 'interstitial', 'fibrosis', 'non-specific',  
176 'pneumonia', 'interstitial fibrosis', 'pneumonitis' and 'fibroelastosis' (Fig. 3m). Further  
177 illustrating the relationship between the age and predicted age of records are kernel  
178 density estimation plots corresponding to each of the topic-specific regressions (Fig.  
179 3n; Extended Data Fig. 3h). Notably, the pulmonary aging topic regression shows  
180 strong age dependency. In sum, our approach effectively leads us to identify an  
181 associated collection of aging modifiers.

## 182 **Cross-validation with PubMed identifies age-modifying drugs**

183 Since we had identified terms in the pathology register that are age-associated we  
184 could identify any other terms (terms, genes, drugs etc.) in other text-based  
185 databases (e.g. PubMed, OMIM.org etc.) that co-occur with these age-associated  
186 pathology terms. We decided to investigate molecules that are co-mentioned in  
187 PubMed abstracts with clinical terms from our identified lung-aging topic (Fig. 4a).  
188 Approximately 35 million molecules in the PubChem library were mined in over 31.8  
189 million PubMed abstracts and assigned a proximity score (Fig. 4b). Among the terms  
190 scoring highest we identified nintedanib, a tyrosine kinase inhibitor, as a potential  
191 pharmacological intervention in aging. Nintedanib is an anti-fibrotic drug used in the  
192



**Fig. 4 | Nintedanib reduces cellular senescence and extends the lifespan of fruit flies.** **a**, Lung aging terms from pathology register combined with molecules in PubMed abstracts. **b**, Proximity scoring of candidate compounds. **c**, Predicted probability of senescence IR and non-IR plates ( $n=3$ , mean mean  $\pm$  SEM). **d**, Percentage of IR exposed cells predicted to be senescent above 0.5 threshold ( $n=3$ , mean mean  $\pm$  SEM). **e**, Relative cell count (normalized to no drug treatment control) in both ionizing radiation (IR) and non-IR plates ( $n=3$ , mean mean  $\pm$  SEM). **f**, RNA seq volcano plots of respective enrichment analyses ( $n=3$ ). **g**, Venn diagram showing common significantly enriched pathways (GSEA) at FDR<0.05 confidence between each case and the respective control DMSO. **h**, EnrichmentMap showing clusters of significantly enriched pathways (GSEA). **i**, *Drosophila melanogaster* survival curves ( $n=90$ ).

194 treatment of idiopathic pulmonary fibrosis<sup>12</sup>. Alongside nintedanib we tested axitinib,  
195 another tyrosine kinase inhibitor with potential anti-fibrotic effects<sup>13</sup>.

## 196 Nintedanib reduces cellular senescence and extends the 197 lifespan of fruit flies

198 To explore whether nintedanib could impact aging we tested the effect of the drug on  
199 cellular senescence, a cellular model of aging that has been implicated in lung  
200 fibrosis<sup>14</sup>. We induced senescence in human dermal fibroblasts by ionizing radiation  
201 (IR) exposure and used our recently published senescence predictor<sup>15</sup> to explore the  
202 effect on the cells. Interestingly, 10  $\mu$ M dose of nintedanib reduced predicted  
203 senescence in IR induced senescent fibroblasts (Fig. 4c,d). However, we also  
204 observed a cytotoxic effect with a 10  $\mu$ M dose of nintedanib in both IR and non-IR  
205 exposed cells manifested in a significant decrease in the relative cell count (Fig. 4e).

206

207 To further understand the effect of nintedanib on senescent cells, we explored  
208 changes in global gene expression through RNA-seq (Fig. 4f). Since nintedanib and  
209 axitinib share common targets<sup>16</sup>, we isolated pathways (Fig. 4g) which were  
210 changed only in senescent cells treated with nintedanib. Notably, nintedanib  
211 downregulated (Fig. 4h) collagen metabolic processes and wound healing gene  
212 pathways which have both been implicated in lung fibrosis and aging<sup>17,18</sup>.

213

214 To explore whether nintedanib could impact aging *in vivo* we investigated the effect  
215 of the drug on the life- and health span of the common aging model organism  
216 *Drosophila melanogaster*, specifically the wild-type *w<sup>1118</sup>* fly. We observed a  
217 significant increase in the maximum lifespan of fruit flies fed a diet containing 100  $\mu$ M

218 dose of nintedanib compared to dimethyl sulfoxide (DMSO) vehicle (Fig. 4i). Notably,  
219 the increase in lifespan is observed in late life. In total, the human pathome allowed  
220 us to identify a drug that may affect the aging process.

## 221 Discussion

222 In this paper we present the human aging pathome (pathoage.com), a compendium  
223 of tissue-specific age- and mortality-associated clinical features extracted from the  
224 clinical text narratives in The Danish Pathology Register. Our investigation shows  
225 strong age-related variance along two trajectories fitting with the ages of  
226 development<sup>19</sup> and with pathological aging<sup>3</sup>. Strikingly, we observed sex-specific  
227 differences in the onset and rate of aging related changes. In males, we observed an  
228 onset of aging related changes around forty years of age, while in females, we  
229 observe that aging trajectories occur almost immediately after development. This  
230 could be considered as evidence towards the hypothesis that aging can be a  
231 selected trait in evolution since it occurs in women prior to peak fertility. Although  
232 speculative, these findings could also suggest that evolution may have allowed  
233 successful males to age later perhaps allowing greater reproduction. It is notable,  
234 that patterns of aging in males occur around the time of mean life-expectancy of  
235 ancient man<sup>20</sup>.

236

237 We assessed whether age could be predicted from clinical text features in lung  
238 records and found relatively poor predictive power considering individual features.  
239 This is not entirely surprising given the abstract nature of language. Nonetheless,  
240 even relatively poor predictive power can reveal useful patterns with the terms  
241 'carcinoma'<sup>21</sup> and 'sarcoidosis'<sup>22</sup> ranked among the most important to prediction



242 accuracy. To improve accuracy, we investigated whether a collection of associated  
243 terms (topics) could contribute to greater predictive power. Indeed, the predicative  
244 power of the topics was approximately 2-fold greater than that of any individual term  
245 and pathology records closely associated with the lung aging topic showed strong  
246 predicative power.

247

248 As an example of the utility of the Pathome, we mined PubMed abstracts for  
249 molecules that occur frequently together with aging lung terms from the pathology  
250 register and identified nintedanib as a potential drug affecting aging. Our  
251 investigation of global gene changes shows that nintedanib down-regulates collagen  
252 metabolism and wound healing pathways in senescent human dermal fibroblasts.  
253 This is compatible with evidence that idiopathic pulmonary fibrosis is characterized  
254 by the accumulation of collagen<sup>18</sup> and an altered wound healing in response to  
255 persistent lung injury<sup>23</sup>. It is important to highlight that the method used to identify  
256 nintedanib can be used to identify any term associated with aging such as the  
257 discovery of new genetic components of aging. The method can also be applied to  
258 identify concepts associated with any pathology described in the database. For  
259 instance, drugs that may impact liver fibrosis, neurodegeneration or any other  
260 defined pathology can be explored.

261

262 In sum, our investigation revealed population-level patterns of aging that are  
263 connected with developmental and pathological aging. This allows us to identify  
264 modifiers of aging that can be translated into new aging interventions. Lastly, we  
265 present The Human Pathome, a unique compendium of thousands of tissue-specific  
266 aging and mortality associated features.

267

## 268 Methods

### 269 Danish dictionary of clinical terms

270 To help identify clinical features in the pathology register we constructed a dictionary  
271 of clinical terms in Danish from the patoSnoMed ontology ([www.patobank/snomed](http://www.patobank/snomed))  
272 and the Danish version of the Systematized Nomenclature of Medicine — Clinical  
273 Terms (SNOMED CT)<sup>24</sup> ontology (<https://sundhedsdatastyrelsen.dk/snomedct>).  
274 Terms found in these ontologies may consist of several words (ex. ‘severe  
275 inflammation’). In addition to using such multi-word terms, we added individual words  
276 from multi word terms to our dictionary (ex. ‘severe’ and ‘inflammation’.)

### 277 Clinical term extraction

278 We identified a total of 2,665,283 unique terms, 178,226 unigrams (one word) and  
279 2,487,957 bigrams (two consecutive words) in 32,961,459 pathology text records in  
280 The Danish National Pathology Register. This yielded a binary matrix of 32,961,459  
281 samples and 2,665,283 features. We filtered this initial dataset keeping only terms  
282 that exist in our dictionary of Danish clinical terms, reducing the size of the dataset to  
283 20,316,270 records and 16,237 terms. We kept terms that appeared at least 50  
284 times in the entire pathology register, and records with 5 or more features present.  
285 We then created individual datasets for tissue specific records. We identified records  
286 associated with specific tissues using a topology (T) code assigned to each record in  
287 the register. For example, to construct a dataset of skeletal system tissues (T10000)  
288 we collected all tissues assigned with a topology code that begins with ‘T1’ thereby  
289 also including bone tissue (T11000). We applied the same filtering strategy used in  
290 the entire dataset to tissue specific datasets. We extracted 242,284 records and

291 4,684 terms for skeletal tissue (T10000), 177,795 records and 4,275 terms for lung  
292 (T28000) and 156,057 records and 4,048 terms for liver (T56000) among other  
293 tissues.

## 294 Term normalization

295 We normalized the clinical term matrix to a term frequency–inverse document  
296 frequency (tf-idf) representation. The tf-idf representation for a term  $t$  in a document  
297  $d$  in a document set consisting of  $n$  documents is  $\text{tf-idf}(t,d)=\text{tf}(t,d)*\text{idf}(t)$ ,  $\text{tf}(t,d)$  being  
298 the frequency of a term  $t$  in document  $d$ , and  $\text{idf}$  being  $\text{idf}(t)=\log[n/\text{df}(t)]+1$  ( $\text{df}(t)$  is  
299 the frequency of term  $t$  in all documents in the document set). To identify the average  
300 term frequency within each age-group we calculated the mean value of all record  
301 vectors within each age group. This yielded one term vector per age group. We  
302 calculated the mean incidence age of clinical terms in the entire register and in each  
303 of the tissue specific datasets.

## 304 Topic modeling with Latent Dirichlet allocation (LDA)

305 We used the scikit-learn implementation of Latent Dirichlet allocation (LDA)<sup>9</sup> to  
306 identify latent semantic structures in the entire corpus of records within tissue  
307 specific datasets. We ran LDA using the batch variational Bayes method. To  
308 determine the optimal number of topics yielding the best fit for the model we  
309 employed a perplexity minimization approach<sup>25</sup> by repeatedly fitting an LDA model to  
310 our dataset and varying the number of topics (2-140). We then identified the number  
311 of topics associated with the smallest perplexity score to be optimal. The topic model  
312 yields topic word distributions signifying the number of times each word is assigned

313 to a topic. Similarly, the topic model also yields document topic distribution signifying  
314 the degree to which a topic is associated with a document.

### 315 Age-prediction from clinical features

316 We used a deep neural network (DNN) multi-layer perceptron (MLP) regression to  
317 predict age from clinical text features in the one-hot representation of the clinical  
318 term matrix. We then calculated the model coefficient of determination score ( $R^2$ )  
319 and the median absolute error (MAE) for each model. We used ordinary least  
320 squares (OLS) linear regression to regress the predicted age and chronological age  
321 of records and calculated topic specific regression slopes. We used the scikit-learn  
322 permutation\_importance function to inspect our DNN age-prediction model to assess  
323 the impact of individual features on the model's accuracy measured by the model's  
324 coefficient of determination score ( $R^2$ ).

### 325 Permutation topic importance

326 To assess the impact of a collection of associated terms (topic) on the model's age-  
327 prediction accuracy we shuffled the collected term vectors within a topic and  
328 calculated the change in the model coefficient of determination score ( $R^2$ ). Topics  
329 associated with a greater change are deemed more important to age-prediction.

### 330 Dimensionality reduction

331 We used the scikit-learn implementation of PCA and t-SNE. We applied PCA and t-  
332 SNE to age-aggregated tf-idf term matrices. We used the python umap-learn  
333 package implementation of UMAP on tf-idf term matrices.

### 334 Term enrichment in tissue and morphology-specific records

335 Term enrichment in tissue or morphology specific records is calculated as  
336  $\log((B+1)/(A-B+1))$  where B is the frequency of a term in tissue or morphology  
337 specific records and A is the frequency a term in the entire dataset.

### 338 PubMed term proximity score

339 We extracted a total of 175,555 unique terms from 31,850,051 PubMed abstracts.  
340 Given a binary feature matrix M and a set A of terms within matrix M we calculated a  
341 proximity score for each individual term in a given set B within matrix M. We applied  
342 a tf-idf transformation to the feature matrix M and calculated the cosine distances  
343 between individual terms in set A to individual terms in set B yielding a distance  
344 matrix AxB. We calculated the term proximity score to be the mean distance of each  
345 term b in set B to all terms in set A that are co-mentioned with term b at least once.

346

347 Matrix M: PubMed abstracts years 2000 onwards.

348 Set A: Aging lung terms.

349 Set B: All PubChem compounds that occur in PubMed abstracts ten times or more.

### 350 Term and topic associated mortality

351 For term and topic associated mortality, we used the R survival package to perform  
352 Cox survival regression. For each clinical term we calculated a Cox regression  
353 coefficient reflecting the hazard associated with the incidence of the term in  
354 pathology records, adjusted for word count and birth year cohort. For topic-  
355 associated mortality, we stratified patient pathology records according to topics. We  
356 created a Boolean variable for each topic reflecting the association of a pathology

357 record with a topic. This yielded a matrix of records and topics. We performed Cox  
358 survival regression on the time to death from examination and noted the Cox  
359 regression coefficient associated with each topic reflecting the hazard associated  
360 with the incidence of the topic in pathology records.

## 361 Cell culture

362 Human primary fibroblast cell lines (Coriell, NJ, USA) AG08498 (AG), GM22159  
363 (159) and GM22222 (222) were cultured in 4.5g/L-enriched Dulbecco's Modified  
364 Eagle's Medium (DMEM)/ Ham's F-12 Nutrient Mix (F12) in a 1:1 solution  
365 supplemented with 10% fetal bovine serum (FBS) and 1% penicillin/streptomycin.  
366 Cells were maintained at 37°C in 5% CO<sub>2</sub> atmosphere conditions and passaged  
367 every 2-3 days. For senescence assays, cells at 70-80% confluency and below 20  
368 passages were seeded in 96-well plates (Corning, 3340) at a density of 3000  
369 cells/well and incubated overnight at 37 °C and 5% CO<sub>2</sub>. Control plates were seeded  
370 at 3000 cells/well or 1500 cells/well. One day after seeding, plates were irradiated  
371 using a YXLON Smart Maxi Shot. Cells were exposed with emission of 0.85 Gy/min.  
372 for 12 minutes for a total exposure of 10 Gy. After IR exposure, cells were incubated  
373 for 6 days with medium changed every 48h. Control plates were seeded on day 6.  
374 On day 7 cells were treated with compounds or vehicle for 48 hours after which the  
375 cells were either harvested for RNA or fixed with 4% paraformaldehyde for 10 min,  
376 washed in PBS and stained with DAPI. Cells were subsequently imaged using an IN  
377 Cell analyzer 2200 high content microscopy at 20x magnification, 12 fields per well.



## 378 RNA sequencing

379 RNA was extracted using Trizol according to manufacturer's protocol. DNBSEQ  
380 Eukaryotic Long Non-Coding RNA-sequencing was performed by BGI Denmark.  
381 Mapping-based quantification of the GRCh38 transcriptome from RNA sequencing  
382 paired-end reads was performed with salmon<sup>26</sup> using a pre-computed transcriptome  
383 index for salmon obtained from refgenie<sup>27</sup>. Differential expression analysis was  
384 performed with DESeq2 1.38.2<sup>28</sup> on genes mapped from transcripts with the  
385 gencode annotation of the Ensembl gene set downloaded from refgenie  
386 [http://refgenomes.databio.org/v3/assets/splash/2230c535660fb4774114bfa966a62f8](http://refgenomes.databio.org/v3/assets/splash/2230c535660fb4774114bfa966a62f823fdb6d21acf138d4/salmon_sa_index?tag=default)  
387 [23fdb6d21acf138d4/salmon\\_sa\\_index?tag=default](http://refgenomes.databio.org/v3/assets/splash/2230c535660fb4774114bfa966a62f823fdb6d21acf138d4/salmon_sa_index?tag=default). Genes with fewer than 10 reads  
388 across all samples were filtered prior to all downstream analyses. Gene set  
389 enrichment analysis was performed using GSEA 4.3.2<sup>29</sup>. An expression dataset file  
390 (.gct) was prepared using DESeq2<sup>28</sup> normalized counts for all samples. Phenotype  
391 labels files (.cls) were prepared for each of the group comparisons. GSEA was run  
392 on with the gene set database 'MSigDB c5.go.bp.v2022.1.Hs.symbols.gmt'<sup>30</sup> and  
393 gene\_set permutation type. We used gene sets which are significantly enriched  
394 (upregulated) at FDR<0.05 in each phenotype in all downstream analyses.  
395 Significantly enriched pathways from GSEA<sup>29</sup> were visualized in Cytoscape<sup>31</sup> using  
396 the EnrichmentMap, AutoAnnotate, WordCloud and clusterMaker2 applications.

## 397 Fruit fly maintenance

398 All diets were made on a standard diet (SD) base consisting of 47.5 g cornflour, 41.6  
399 g dextrose, 19.3 g Brewer's Yeast, 6.55 g Low Melting Agar (Calbiochem), and  
400 2.46% Nipagin (Merck, Germany) per litre. All ingredients except Nipagin were mixed  
401 and heated to 80°C. When the mixture had cooled to 40°C, Nipagin was added. The

402 mix was distributed in falcon tubes and compounds added in various concentrations  
403 to make the treatment diets. Diets with equivalent amounts of DMSO were used as  
404 controls. Stock flies were housed in vials of 30 flies to avoid overcrowding and kept  
405 on the standard diet. Both stock and treatment flies were kept at a constant  
406 temperature of 25°C, a relative humidity of 60%, and a 12:12 h light:dark cycle. The  
407 wild type strain *w<sup>1118</sup>* (Bloomington Drosophila Stock Center) was used for all  
408 longevity assay. Before assays, 5-10 crosses with a ratio of 15:9 female to male flies  
409 were set and kept under standard rearing conditions in polypropylene vials in  
410 standard diet.

#### 411 Fruit fly lifespan assay

412 Every three days, for 9-12 days, flies were flipped into new vials containing the  
413 standard diet. Hatches were collected at birth and put in new vials with the desired  
414 compound condition. Per each condition, three vials with ten male flies each were  
415 prepared. Vials were put in front of cameras for our fly tracking system as part of the  
416 Tracked.bio platform ([www.tracked.bio](http://www.tracked.bio)). For all longevity assay, flies were flipped  
417 once weekly into vials with freshly prepared food. Each vial had ten male flies, which  
418 were selected from the new-born hatches of the set crosses. Male flies were chosen  
419 among those which did not show any damage to the wings.  
420 During each flipping, flies were counted, and data collected into a spreadsheet.  
421 Behavioral metrics were calculated from the Tracked.bio system. We used the  
422 lifelines python package to fit a Kaplan-Meier estimator for the survival function of  
423 fruit fly lifespan and to perform a log-rank test to test for statistically significant  
424 differences in survival. A count of live flies in vials was recorded once per week.  
425 Since fruit fly vials were initiated over a period of several days as newly hatched flies

426 were collected, we extrapolated weekly counts to daily counts before performing  
427 survival analysis.

428

429

## 430 Acknowledgements

431 This research was supported by the Novo Nordisk Foundation Challenge  
432 Programme (#NNF17OC0027812), the Nordea Foundation (#02-2017-1749), the  
433 Neye Foundation, the Lundbeck Foundation (#R324-2019-1492), the Ministry of  
434 Higher Education and Science (#0238-00003B) and Insilico Medicine. The funders  
435 had no role in study design, data collection and analysis, decision to publish or  
436 preparation of the manuscript.

437

438

## 439 References

- 440 1. López-Otín, C., Blasco, M. A., Partridge, L., Serrano, M. & Kroemer, G. Hallmarks of  
441 aging: An expanding universe. *Cell* **186**, 243–278 (2023).
- 442 2. Andreassen, S. N., Michael Ben Ezra & Scheibye-Knudsen, M. A defined human aging  
443 phenome. *Aging* **11**, 5786–5806 (2019).
- 444 3. Niccoli, T. & Partridge, L. Ageing as a risk factor for disease. *Curr. Biol. CB* **22**, R741-752  
445 (2012).
- 446 4. Névéol, A., Dalianis, H., Velupillai, S., Savova, G. & Zweigenbaum, P. Clinical Natural  
447 Language Processing in languages other than English: opportunities and challenges. *J.*  
448 *Biomed. Semant.* **9**, 12 (2018).
- 449 5. Jensen, P. B., Jensen, L. J. & Brunak, S. Mining electronic health records: towards better  
450 research applications and clinical care. *Nat. Rev. Genet.* **13**, 395–405 (2012).
- 451 6. Denny, J. C. *et al.* Systematic comparison of phenome-wide association study of  
452 electronic medical record data and genome-wide association study data. *Nat. Biotechnol.*  
453 **31**, 1102–1110 (2013).
- 454 7. Roque, F. S. *et al.* Using electronic patient records to discover disease correlations and  
455 stratify patient cohorts. *PLoS Comput. Biol.* **7**, e1002141 (2011).
- 456 8. Erichsen, R. *et al.* Existing data sources for clinical epidemiology: the Danish National  
457 Pathology Registry and Data Bank. *Clin. Epidemiol.* **2**, 51–56 (2010).
- 458 9. Blei, D. M., Ng, A. Y. & Jordan, M. I. Latent Dirichlet Allocation. *J. Mach. Learn. Res.* **3**,  
459 993–1022 (2003).
- 460 10. Wollin, L., Maillet, I., Quesniaux, V., Holweg, A. & Ryffel, B. Antifibrotic and anti-  
461 inflammatory activity of the tyrosine kinase inhibitor nintedanib in experimental models of  
462 lung fibrosis. *J. Pharmacol. Exp. Ther.* **349**, 209–220 (2014).
- 463 11. Rogers, L. K. & Cismowski, M. J. Oxidative Stress in the Lung - The Essential  
464 Paradox. *Curr. Opin. Toxicol.* **7**, 37–43 (2018).

- 465 12. Richeldi, L. *et al.* Efficacy and safety of nintedanib in idiopathic pulmonary fibrosis. *N.*  
466 *Engl. J. Med.* **370**, 2071–2082 (2014).
- 467 13. Richeldi, L. *et al.* Efficacy of a tyrosine kinase inhibitor in idiopathic pulmonary  
468 fibrosis. *N. Engl. J. Med.* **365**, 1079–1087 (2011).
- 469 14. Schafer, M. J. *et al.* Cellular senescence mediates fibrotic pulmonary disease. *Nat.*  
470 *Commun.* **8**, 14532 (2017).
- 471 15. Heckenbach, I. *et al.* Nuclear morphology is a deep learning biomarker of cellular  
472 senescence. *Nat. Aging* **2**, 742–755 (2022).
- 473 16. Slobbe, P. *et al.* Two anti-angiogenic TKI-PET tracers, [(11)C]axitinib and  
474 [(11)C]nintedanib: Radiosynthesis, in vivo metabolism and initial biodistribution studies in  
475 rodents. *Nucl. Med. Biol.* **43**, 612–624 (2016).
- 476 17. Maher, T. M., Wells, A. U. & Laurent, G. J. Idiopathic pulmonary fibrosis: multiple  
477 causes and multiple mechanisms? *Eur. Respir. J.* **30**, 835–839 (2007).
- 478 18. Jessen, H. *et al.* Turnover of type I and III collagen predicts progression of idiopathic  
479 pulmonary fibrosis. *Respir. Res.* **22**, 205 (2021).
- 480 19. Coleman, L. & Coleman, J. The measurement of puberty: a review. *J. Adolesc.* **25**,  
481 535–550 (2002).
- 482 20. Eshed, V., Gopher, A., Gage, T. B. & Hershkovitz, I. Has the transition to agriculture  
483 reshaped the demographic structure of prehistoric populations? New evidence from the  
484 Levant. *Am. J. Phys. Anthropol.* **124**, 315–329 (2004).
- 485 21. Torre, L. A., Siegel, R. L. & Jemal, A. Lung Cancer Statistics. *Adv. Exp. Med. Biol.*  
486 **893**, 1–19 (2016).
- 487 22. Varron, L., Cottin, V., Schott, A.-M., Broussolle, C. & Sève, P. Late-onset sarcoidosis:  
488 a comparative study. *Medicine (Baltimore)* **91**, 137–143 (2012).
- 489 23. Zhang, L. *et al.* Macrophages: friend or foe in idiopathic pulmonary fibrosis? *Respir.*  
490 *Res.* **19**, 170 (2018).
- 491 24. Lee, D., de Keizer, N., Lau, F. & Cornet, R. Literature review of SNOMED CT use. *J.*  
492 *Am. Med. Inform. Assoc. JAMIA* **21**, e11-19 (2014).

- 493 25. Zhao, W. *et al.* A heuristic approach to determine an appropriate number of topics in  
494 topic modeling. *BMC Bioinformatics* **16 Suppl 13**, S8 (2015).
- 495 26. Patro, R., Duggal, G., Love, M. I., Irizarry, R. A. & Kingsford, C. Salmon provides fast  
496 and bias-aware quantification of transcript expression. *Nat. Methods* **14**, 417–419 (2017).
- 497 27. Stolarczyk, M., Reuter, V. P., Smith, J. P., Magee, N. E. & Sheffield, N. C. Refgenie:  
498 a reference genome resource manager. *GigaScience* **9**, giz149 (2020).
- 499 28. Love, M. I., Huber, W. & Anders, S. Moderated estimation of fold change and  
500 dispersion for RNA-seq data with DESeq2. *Genome Biol.* **15**, 550 (2014).
- 501 29. Subramanian, A. *et al.* Gene set enrichment analysis: A knowledge-based approach  
502 for interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci.* **102**, 15545–  
503 15550 (2005).
- 504 30. Liberzon, A. *et al.* The Molecular Signatures Database (MSigDB) hallmark gene set  
505 collection. *Cell Syst.* **1**, 417–425 (2015).
- 506 31. Shannon, P. *et al.* Cytoscape: a software environment for integrated models of  
507 biomolecular interaction networks. *Genome Res.* **13**, 2498–2504 (2003).
- 508
- 509



## 510 Figure Legends

511

512 **Fig. 1 | Pathological aging begins post-development for females and at mid-life**

513 **for males. a**, The Human Aging Pathome and aging intervention discovery concept

514 and workflow. **b**, PCA of age-aggregated pathology records (n=20,316,270) from in

515 entire pathology register. Normalized Euclidean distance between age adjacent PCA

516 coordinates. **c,d**, PC1, PC2 coordinates vs. Age. **e**, PCA of age-aggregated of

517 pathology records (n=14,492,989) from females in the entire pathology register.

518 Normalized Euclidean distance between age adjacent PCA coordinates. **f,g**, PC1,

519 PC2 coordinates vs. Age. **h**, PCA of age-aggregated of pathology records

520 (n=5,823,281) from males in the entire pathology register. Normalized Euclidean

521 distance between age adjacent PCA coordinates. **i,j**, PC1, PC2 coordinates vs. Age.

522 **k**, t-SNE of clinical features in age-aggregated pathology records in the entire

523 pathology register. **l**, UMAP of clinical features in pathology records in the entire

524 pathology register. **m-o**, Positive morphology specific enrichment: Inflammation,

525 Adenocarcinoma and Adenoma. **p-r**, Positive enrichment of clinical terms in tissue

526 specific pathology records: Lung, Skeletal system and Bone marrow.

527

528 **Fig. 2 | Tissues age along specific trajectories. a**, PCA of age-aggregated tissues

529 specific pathology records. **b**, UMAP of clinical features of tissue specific pathology

530 records (mean incidence age). **c**, UMAP of clinical features of tissue specific

531 pathology records (Cox regression coefficient). Tissues shown are: Lung

532 (n=177,795), Liver (n=156,057), Heart (n=27,055), Kidney (n=85,244), Nervous

533 system (n=183,729), Bladder (n=250,532), Skeletal system (n=242,282) and

534 Gallbladder (n=182,261). **d**, Normalized Euclidean distance between age-adjacent  
535 PCA coordinates of all tissues in. **e-h**, Positive enrichment of clinical terms in  
536 morphology specific lung records: **e**, Inflammation. **f**, Benign tumor. **g**, Adenoma. **h**,  
537 Adenocarcinoma.

538

539 **Fig. 3 | Semantic structures in clinical text describe lung pathologies and**  
540 **predict age. a**, DNN age-prediction from clinical features in lung pathology records.  
541 **b**, Permutation feature importance (PFI). **c**, UMAP of clinical features in lung records  
542 (LDA topic). **d**, t-SNE of clinical feature LDA distributions in topics (LDA topics). **e**, t-  
543 SNE of LDA record distributions in topics (LDA topics). **f**, t-SNE of clinical feature  
544 LDA distributions in topics (mean incidence age). **g**, t-SNE of clinical feature LDA  
545 distributions in topics (mortality: Cox regression coefficient). **h**, t-SNE of LDA record  
546 distributions in topics (age). **i**, Topic associated age-prediction linear regression. **j**,  
547 Topic associated permutation feature importance (PFI). **k**, Topic associated age-  
548 prediction regression slope. **l**, Terms describing Obstetric Gynecologic Pulmonary  
549 Complications (OGPC). **m**, Terms describing pulmonary aging (PA). **n**, Bivariate  
550 kernel density estimation (kde) and histograms of chronological age vs. predicted  
551 age, records closely associated with PF and OGPC.

552

553 **Fig. 4 | Nintedanib reduces senescence in cells in culture and extends the**  
554 **lifespan of fruit flies. a**, Lung aging terms from pathology register combined with  
555 molecules in PubMed abstracts. **b**, Proximity scoring of candidate compounds. **c**,  
556 Relative cell count (normalized to no drug treatment control) in both ionizing radiation  
557 (IR) and non-IR plates (n=3, mean mean  $\pm$  SEM). **d**, Percentage of IR exposed cells  
558 predicted to be senescent above 0.5 threshold (n=3, mean mean  $\pm$  SEM). **e**,

559 Predicted (deep neural network) probability of senescence IR and non-IR plates  
560 (n=3, mean mean  $\pm$  SEM). **f**, RNA seq volcano plots of respective enrichment  
561 analyses (n=3). **g**, Venn diagram showing common significantly enriched pathways  
562 (GSEA) at FDR<0.05 confidence between each case and the respective control  
563 DMSO. **h**, EnrichmentMap showing clusters of significantly enriched pathways  
564 (GSEA). **i**, *Drosophila melanogaster* survival curves (n=90).

565

566 **Extended Data Fig. 1 | Sex-specific patterns.** **a**, UMAP of clinical features in  
567 pathology records from males in the entire pathology register. **b**, UMAP of clinical  
568 features in pathology records from females in the entire pathology register. **c**, t-SNE  
569 of clinical features in age-aggregated pathology records from males in the entire  
570 pathology register. **d**, t-SNE of clinical features in age-aggregated pathology records  
571 from females in the entire pathology register. **e**, Positive enrichment of clinical terms  
572 in various tissue and morphology specific records **f**, Term count associated hazard.  
573 **g**, Birth cohort associated hazard.

574

575 **Extended Data Fig. 2 | Tissue-specific analyses.** **a**, PCA of age-aggregated  
576 tissues specific pathology records. **b**, UMAP of clinical features of tissue specific  
577 pathology records (mean incidence age). **c**, UMAP of clinical features of tissue  
578 specific pathology records (Cox regression coefficient). **d**, Normalized Euclidean  
579 distance between age adjacent PCA coordinates of all tissues in.

580

581 **Extended Data Fig. 3 | Topic modelling.** **a**, LDA perplexity for model fitted with  
582 varying number of topics (skeletal system, lung, liver). **b**, Topic associated mean  
583 age. **c**, Topic associated mortality. **d**, Topic associated age-prediction linear

584 regression. **e**, Topic associated permutation feature importance (PFI). **f**, Topic  
585 associated age-prediction regression slope. **g**, Terms describing Human  
586 Immunodeficiency Virus (HIV). **h**, Bivariate kernel density estimation (kde) and  
587 histograms of chronological age vs. predicted age, records closely associated with  
588 HIV and t-SNE of clinical feature LDA distributions in topics highlighting OGPC, PA  
589 and HIV topics.