

Supplementary Discussion, Methods, Figures, Tables and References

Table of contents

Table of contents	1
Results	3
Genome Wide Association Studies	3
Musculoskeletal general aging	3
Full body aging	3
Spine aging	4
Hip aging	4
Knee aging	5
Anthropometry, impedance, heel bone densitometry and hand grip strength-based aging	5
Discussion	6
Comparison between our age predictors and the literature	6
Full body X-rays	6
Hip X-rays	7
Knee X-rays	8
Hand X-rays	8
Other skeletal structures	9
Methods	9
Hardware	9
Software	9
Training, tuning and predictions	10
Data splitting	10
Scalar data	11
Nested cross-validation	11
Bayesian hyperparameters optimization	11
Example	12
Images	15
Hyperparameters tuning upstream of the cross-validation	15
Cross-validation	18

Cross-validation example	18
Interpretability of the predictions	19
Scalar data-based predictors	19
Image-based predictors	20
Non-genetic correlates of accelerated aging	23
Imputation of the non-genetic X-variables	24
X-Wide Association Studies	25
Supplementary Figures	27
Supplementary Tables	41
Supplementary References	50

Results

Genome Wide Association Studies

Musculoskeletal general aging

Accelerated musculoskeletal aging is $34.9 \pm 1.8\%$ heritable and 11 SNPs in nine genes are significantly associated with this phenotype (Figure 5). The GWAS highlighted six peaks: (1) FGFR3P1 (Fibroblast Growth Factor Receptor 3 Pseudogene 1, a pseudogene involved in vascular endothelial growth factor signaling); (2) TBX15 (T-Box Transcription Factor 15, transcription factor involved in development and linked to Cousin syndrome and acromegaloid facial appearance syndrome); (3) CCDC91 (Coiled-Coil Domain Containing 91, involved in the trans-Golgi network associated with ossification of the posterior longitudinal ligament of the spine and diffuse idiopathic skeletal hyperostosis); (4) AAGAB (Alpha And Gamma Adaptin Binding Protein, involved in vesicle trafficking and linked to palmoplantar keratoderma); (5) FAM3C (FAM3 Metabolism Regulating Signaling Molecule C, linked to deafness and pancreatic cancer); and (6) SCEL (Sciellin, expressed in keratinocytes and linked to extrahepatic bile duct adenocarcinoma).

Full body aging

Accelerated full body aging is $30.7 \pm 1.6\%$ heritable and eight SNPs in eight genes are significantly associated with this phenotype (Fig. S6). The GWAS highlighted four peaks: (1) TBX15 (T-Box Transcription Factor 15, transcription factor involved in development and linked to Cousin syndrome and acromegaloid facial appearance syndrome); (2) FGFR3P1 (Fibroblast

Growth Factor Receptor 3 Pseudogene 1, a pseudogene involved in vascular endothelial growth factor signaling); (3) SCEL (Sciellin, expressed in keratinocytes and linked to extrahepatic bile duct adenocarcinoma); and (4) FGFR1 (Fibroblast Growth Factor Receptor 1, a fibroblast growth factor receptor linked to osteoglophonic dysplasia and Hartsfield syndrome).

Spine aging

Accelerated spine aging is $32.9 \pm 1.7\%$ heritable and 22 SNPs in 13 genes are significantly associated with this phenotype (Fig. S7). The GWAS highlighted nine peaks: (1) TNFSF11 (Tumor Necrosis Factor Ligand Superfamily Member 11, a ligand for osteoprotegerin involved in osteoclast differentiation and activation, and linked to osteopetrosis); (2) HTRA1 (HtrA Serine Peptidase 1, involved in cell growth regulation); (3) FAM3C (FAM3 Metabolism Regulating Signaling Molecule C, linked to deafness and pancreatic cancer); (4) SUPT3H (SPT3 Homolog, SAGA And STAGA Complex Component, involved in chromatin organization and linked to dysostosis and hepatic adenoma); (5) RP11-438P9.2 (a pseudogene); (6) SOX5 (SRY-Box Transcription Factor 5, involved in chondrocytes differentiation and cartilage formation, linked to Lamb-Shaffer syndrome); (7) ESR1 (Estrogen Receptor 1, an estrogen receptor involved in sexual functions and linked to osteoporosis, breast cancer and endometrial cancer); (8) COLGALT2 (Collagen Beta(1-O)Galactosyltransferase 2, involved in degradation of the extracellular matrix and linked to porencephaly); and (9) AAGAB (Alpha And Gamma Adaptin Binding Protein, involved in vesicle trafficking and linked to palmoplantar keratoderma).

Hip aging

Accelerated hip aging is $27.7 \pm 1.6\%$ heritable and 20 SNPs in seven genes are significantly associated with this phenotype (Fig. S8). The GWAS highlighted four peaks: (1) SUPT3H (SPT3

Homolog, SAGA And STAGA Complex Component, involved in chromatin organization and linked to dysostosis and hepatic adenoma); (2) CCDC91 (Coiled-Coil Domain Containing 91, involved in the trans-Golgi network associated with ossification of the posterior longitudinal ligament of the spine and diffuse idiopathic skeletal hyperostosis); (3) TNFSF11 (Tumor Necrosis Factor Ligand Superfamily Member 11, a ligand for osteoprotegerin involved in osteoclast differentiation and activation, and linked to osteopetrosis); and (4) RP11-396020.1 (a long intergenic non-coding RNA).

Knee aging

Accelerated knee aging is $25.3 \pm 1.7\%$ heritable and 16 SNPs in six genes are significantly associated with this phenotype (Fig. S9). The GWAS highlighted four peaks: (1) WNT16 (Wnt Family Member 16, involved in signaling, cell fate and oncogenesis); (2) CCDC91 (Coiled-Coil Domain Containing 91, involved in the trans-Golgi network associated with ossification of the posterior longitudinal ligament of the spine and diffuse idiopathic skeletal hyperostosis); and (3) HLA-DRB9 (Major Histocompatibility Complex, Class II, DR Beta 9 Pseudogene, a pseudogene).

Anthropometry, impedance, heel bone densitometry and hand grip strength-based aging

Accelerated musculoskeletal scalar features-based (anthropometry, impedance, heel bone densitometry and hand grip strength) aging is $22.1 \pm 0.2\%$ heritable and that 2,631 SNPs in 954 genes are significantly associated with this phenotype (Fig. S10). The ten highest peaks highlighted by the GWAS are (1) GDF5 (Growth Differentiation Factor 5, involved in the development of cartilage, joints, brown fat, teeth, and neurons, linked to osteoarthritis,

acromesomelic dysplasia, brachydactyly, proximal symphalangism, chondrodysplasia and multiple synostoses syndrome); (2) CRHR1 (Corticotropin Releasing Hormone Receptor 1, a GPCR protein involved in the hypothalamic-pituitary-adrenal pathway); (3) HOTTIP (HOXA Distal Transcript Antisense RNA, involved in cell proliferation and linked to tumors); (4) SOCS2 (Suppressor Of Cytokine Signaling 2, involved in cytokine signaling and linked to polycythemia); (5) PKD2L1 (Polycystin 2 Like 1, Transient Receptor Potential Cation Channel, involved in cell-cell and cell-matrix interactions and linked to polycystic kidney disease); (6) GPR126 (Adhesion G Protein-Coupled Receptor G6, a G protein-coupled receptor involved in body height, linked to distal arthrogyrosis); (7) AC013480.2 (an antisense gene). AC013480.2 is in linkage disequilibrium with PKDCC (Protein Kinase Domain Containing, Cytoplasmic, involved in development and linked to rhizomelic limb shortening with dysmorphic features); (8) POLD3 (DNA Polymerase Delta 3, Accessory Subunit, involved in DNA replication and repair and linked to Ruijs-Aalfs syndrome); (9) ZBTB38 (Zinc Finger And BTB Domain Containing 38, a transcriptional activator involved in height); and (10) FUBP3 (Far Upstream Element Binding Protein 3, involved in gene expression).

Discussion

Comparison between our age predictors and the literature

Full body X-rays

Langner et al. predicted chronological age with a testing R^2 of 82% and MAE of 2.47 ± 1.91 years after training the VGG16 CNN architecture¹ on 32,323 UKB's full body Dual-energy X-ray

absorptiometry [DXA] images². We found that our preprocessing and the architectures that we used did not outperform their model individually. However, we found that our ensemble model yielded a statistically significant higher final testing R^2 of $85.7\pm 0.1\%$ (RMSE= 2.85 ± 0.01 years). We also found that adding to this ensemble predictions generated on spine DXA images ($R^2=74.6\pm 0.2\%$), hip DXA images ($R^2=69.0\pm 0.3\%$), Knee DXA images ($R^2=69.0\pm 0.3\%$) and scalar musculoskeletal features such as anthropometry, heel bone densitometry and hand grip strength ($R^2=25.9\pm 0.1\%$) yielded a chronological age predictor with a R^2 of $87.6\pm 0.1\%$.

Karargyris et al. used chest X-rays and the DenseNet121 CNN architecture to predict chronological age with a R^2 of 90% ³. After testing several architectures on the ChestX-ray8⁴ dataset using transfer learning we found that the best architecture was indeed DenseNet121, which yielded the same R^2 of 90% . We also found that our ensemble model predicted chronological age with a R^2 of 91% . These R^2 values are slightly higher than the one we obtained on the UKB datasets ($85.7\pm 0.1\%$ for the full body DXA images), but as previously explained this difference in prediction accuracy is largely driven by the difference in age ranges. Figure 1 of the publication shows that the age range for the ChestX-ray8 dataset is 1-90 years, which includes the childhood and the teenagerhood during which significant anatomical changes happen, compared to 37-82 years for UKB.

Hip X-rays

Wittschieber et al.⁵ used a linear regression on features extracted from 643 pelvic X-rays collected from participants aged 10-30 years and obtained a R^2 value of 38% . Similar work was performed on the pelvis⁶, the iliac crest^{7,8} and the ilium⁹. Our model ($R^2=69.0\pm 0.3\%$; RMSE= 4.20 ± 0.02) significantly outperformed the one developed by Wittschieber et al. despite

our age range not covering teenagerhood years, which can be explained by the 65 times larger sample size we leveraged.

Knee X-rays

Two chronological age predictors were built on knee X-ray images collected from young participants^{10,11}. O'Connor et al. used a linear regression on features extracted from 221 knee X-rays collected from participants aged 9-19 years and obtained a R^2 value of 77.5-81.5%¹⁰. Tang et al. used a linear regression on features extracted from 503 knee X-rays collected from participants aged 6-19 years and obtained a R^2 value of 88-90%¹¹. Fan et al. used a linear regression on features extracted from 322 knee X-rays and MRIs collected from participants aged 11-30 years and obtained R^2 values ranging from 44.1% to 65.4%¹² and other researchers performed similar studies¹³⁻¹⁶. We trained CNNs on 79,477 knee X-ray images collected from participants aged 45-82 and obtained a R^2 value of $69.0 \pm 0.3\%$ and a RMSE value of 4.20 ± 0.02 years. The significantly higher accuracy of the two first models described above can be explained by the age range over which they were trained and tested, as the human body changes significantly faster during development. We are, to our knowledge, the first to build a chronological age from knee X-rays for the older adult population.

Hand X-rays

Several researchers leveraged hand X-ray images to predict chronological age in infants, children, and teenagers during the last six years¹⁷⁻²⁴. Westerberg for example trained the Xception²⁵ architecture on 12,811 images from patients aged 0-19 years and predicted chronological age with a MAE of 1.007 years¹⁹.

Other skeletal structures

Other skeletal structures have been leveraged to predict age such as the clavicle epiphyses^{26–28}, the fourth cervical vertebra²⁹, the mandibular ramus length³⁰ or the molars^{31–33}.

Methods

Hardware

We performed the computation for this project on Harvard Medical School's compute cluster, with access to both central processing units [CPUs] and general processing units [GPUs] (Tesla-M40, Tesla-K80, Tesla-V100) via a Simple Linux Utility for Resource Management [SLURM] scheduler.

Software

We coded the project in Python³⁴ and used the following libraries: NumPy^{35,36}, Pandas³⁷, Matplotlib³⁸, Plotly³⁹, Python Imaging Library⁴⁰, SciPy^{41–43}, Scikit-learn⁴⁴, LightGBM⁴⁵, XGBoost⁴⁶, Hyperopt⁴⁷, TensorFlow 2⁴⁸, Keras⁴⁹, Keras-vis⁵⁰, iNNvestigate⁵¹. We used Dash⁵² to code the website on which we shared the results. We set the seed for the os library, the numpy library, the random library and the tensorflow library to zero.

Training, tuning and predictions

Data splitting

We split the 676,787 samples into ten data folds, while keeping all samples from the same participant in the same fold. To ensure this, we split the 502,211 participants' ids (referred to by UKB as "eid") into ten different buckets of the same size. To generate ten folds for each sub-dataset (e.g. ECGs), we took the intersection of the samples in each of the ten folds with the samples for which the sub-dataset data was available. This method had however one important loophole, which is that we could not guarantee that the folds for the sub-datasets would be balanced. For example, resting ECG data was only recorded for 42,360 out of the 502,211 participants. Since the 502,211 participants are split into ten folds, a fold contains approximately 50,221 participants. Although unlikely, we could therefore not guarantee that all or most of the ECG samples would be attributed to the first data fold, leading to highly unbalanced folds for the ECGs analysis. Unbalanced folds can lead to problems during the cross-validation (see further below), as models trained on a smaller number of samples will tend to generalize worse. One solution would have been to use a different split for each dataset, but this would have generated problems when building the ensemble models fold by fold (see Methods - Models ensembling). To mitigate this issue of unbalanced data folds, we developed the following heuristic. We randomly split the 502,211 participants into ten folds, 1,000 times. For each of these 1,000 splits, we computed for each sub-dataset the variance of the percentages of samples in each fold. We then scored each of the 1,000 splits using the maximum of the variance among the different sub-datasets. For example, if the ECG samples were not evenly split for the i th split out of the 1,000 splits (e.g. fold 1: 55% of the samples,

every other fold: 5% of the samples), the variance of the sample proportions would be high, which would yield a poor score for the i th split. Finally, we selected the split with the lowest score as the final split for the main dataset, and for all the sub-datasets. This selected split had a score of $5.8e-4$, which means that the most unbalanced sub-dataset had a variance in its sample size proportion between its ten folds of $5.8e-4$.

Scalar data

Nested cross-validation

Cross-validation is a method to tune the regularization of models and prevent overfitting⁵³. For the models inputting scalar data (Figure 1A in green), we tuned the hyperparameters and generated a testing prediction for each sample using a nested 10x9-folds cross-validation. We refer to the two nested cross-validations as the “outer” and the “inner” cross-validations. The outer-cross validation is used to generate an unbiased testing prediction for each sample, as opposed to a simple split of the data into a “training+validation” set on one hand, and a testing set on the other hand, which would only generate a testing prediction for one tenth of the dataset. The inner cross-validation is used to tune the hyperparameters more precisely, leveraging the full inner cross-validation dataset as a validation set, as opposed to a simple data split of the “training+validation” dataset into a training and a validation sets, which would only use one data fold as the validation set to estimate the performance associated with a specific combination of hyperparameters. The nested cross-validation is illustrated in Table S26.

Bayesian hyperparameters optimization

To tune the hyperparameters, we used the Tree-structured Parzen Estimator Approach⁵⁴ [TPE] of the hyperopt python package⁵⁵. TPE is a sequential Bayesian hyperparameters optimization

method that iteratively suggests the next most promising hyperparameters combination as a function of the hyperparameters combinations that have already been tested, by building a probabilistic representation of the objective function. We set the number of iterations to 30. For each model, 30 different hyperparameter combinations are iteratively tested before selecting the best performing one. The hyperparameters names and their ranges defining the hyperparameters space can be found in Table S25. It might be of interest to other researchers that we initially tuned the hyperparameters using a random search⁵⁶ with the same number of iterations, and we did not observe a significant improvement in the model's performance after implementing the Bayesian hyperparameters optimization.

Example

For the sake of clarity, let us walk through a concrete example, which is illustrated in Table S26. Suppose we want to generate unbiased predictions for every sample in a dataset using an elastic net. First, let us generate the testing prediction for the data fold F9, which is performed by the first fold of the outer cross-validation (outer cross-validation fold 0). We select the data fold F9 out of the ten data folds as the testing fold, and we select the remaining nine data folds as "training+validation" folds for the inner cross-validation. We scale and center the target (age) and the predictors using the mean and standard deviation values of the variables on the "training+validation" dataset. We then enter the first inner-cross validation.

For the first inner cross-validation fold, we select the data fold F8 as the validation set, and the remaining eight "training+validation" data folds as the training set. We re-scale and center age and the predictors in the training and the validation sets using the mean and standard deviation values of the training set. We train the model on the eight training data folds with the first

hyperparameters combination sampled by the TPE algorithm (one value for alpha and one value for l1_ratio) and generate validation predictions on the validation fold (data fold F8), which we unscale. This completes the first of the nine inner cross-validation folds (Inner CV fold 0). We then permute the nine inner data folds. We scale the age and the predictors using the mean and standard deviation computed on the new training set. Then we train the model with the same first combination of hyperparameters on eight data folds, leaving aside the data fold F9 (still being used as the testing set for the outer cross-validation) and the data fold F7 (now being used as the validation set for the inner cross-validation). We then use the new trained model to generate validation predictions on the data fold F7, which we unscale. This completes the second of the nine inner-cross validation folds (Inner CV fold 1). We then reiterate these inner permutation and training processes seven more times, until every data fold in the nine “training+validation” data folds is used as the validation set once. At this point, we concatenate the validation predictions from these nine validation folds to obtain the overall validation predictions associated with the first hyperparameters combination, and compute the associated performance metric (e.g. RMSE). This completes the inner-cross validation for the first hyperparameters combination.

We then perform the same 9-folds inner cross-validation, this time with the second hyperparameters combination suggested by the TPE algorithm. We iterate this process 28 more times, until 30 different hyperparameters combinations have iteratively been tested. Next, we select the hyperparameter combination that yielded the best validation performance (e.g. minimum RMSE), and we retrain a model on the whole nine “training+validation” data folds (all data folds except for data fold #1), using this best performing hyperparameters combination. This completes the first inner cross-validation.

We then use the model to generate unbiased predictions on the unseen testing set (data fold F9) and record these predictions. By anticipation for the ensembling algorithm (see Methods - Models ensembling) we also need to compute validation predictions on the data fold F8. We do this by training a model on all the data folds aside from the validation fold (data fold F8) and the testing fold (data fold F9), with the selected hyperparameters combination. We then use this trained model to compute predictions on the validation fold (data fold F8) and record these predictions, after unscaling them. This completes the first of the ten outer cross-validation folds (outer cross-validation 0).

We then complete the second outer cross-validation fold (outer cross-validation 1), this time using the data fold F8 as the testing dataset, to obtain unbiased testing predictions on this data fold, as well as validation predictions on the data fold F7. We reiterate the process eight more times to obtain the testing and validation predictions on the remaining data folds. We then concatenate the testing predictions from the ten data folds to obtain our final testing predictions for the model. Similarly, we concatenate the validation predictions from the ten data folds to obtain our final testing predictions for the model, which will later be used during ensemble models building and model selection (see Methods - Models ensembling).

The final validation and testing predictions for each data fold are therefore not necessarily associated with the same hyperparameters combination. It is also important to notice that we performed a single outer cross-validation, but that we performed a separate inner-cross validation for each outer cross-validation fold (hence the word “nested”), for a total of ten inner cross-validations per outer cross-validation fold.

Images

Hyperparameters tuning upstream of the cross-validation

The hyperparameters we tuned were the number of added fully connected dense layers, the number of nodes in these layers, their activation function, the optimizer, the initial learning rate, the weight decay, the dropout rate, the data augmentation amplitude and the batch size. Repeatedly tuning the values of the hyperparameters for different deep neural networks architectures and on the different cross-validation folds would have been prohibitively time and resource consuming. Instead, we sequentially explored how each hyperparameter was affecting the training and validation performances for a single architecture (InceptionV3) on a single cross validation fold (fold #0, see [Methods - Training, tuning and predictions - Images - Cross-validation](#) for the detailed description of the cross-validation). We then extrapolated the hyperparameter values to the other architectures, datasets and cross-validation folds. The hyperparameters combinations tested during the tuning can be found in [Table S27](#).

First, we maximized the batch size for each architecture. The maximum number of images per batch depends on the memory of the GPU and the size of the architecture, which itself depends on the dimensions of the image. We used a batch size of 32 for InceptionV3 and 8 for InceptionResNetV2.

Then, we tested the learning rates, including $1e-6$, $1e-5$, $1e-4$, $1e-3$, $1e-2$ and $1e-1$. We observed that learning rates larger than $1e-4$ prevented the model from converging for some runs. Second, we did not observe significant differences between the results obtained with learning rates smaller than $1e-4$. We therefore set the initial learning rate to be $1e-4$ for all

models to shorten the time to convergence while ensuring that the learning rate was small enough to allow convergence and the finding of a local minima for the loss function.

Then we tested three different optimizers to perform the gradient descent: Adam ⁵⁷, Adadelta ⁵⁸ and RMSprop ⁵⁹. We did not observe any significant differences between the optimizers, so we set the optimizer to be Adam.

We then added different numbers of fully connected layers between the base CNN and side CNN's concatenated outputs and the final activation layer. We set the number of nodes to be 1,024 in the first added layer and then decreased the number of nodes by a factor of two for each successive layer. For example, if we added three fully connected layers, the number of nodes was 1024, 512 and 256. We added zero, one and five layers. We did not observe significant differences in the performance of the different architectures, so we set the number of fully connected layers to one.

We then tested powers of two from 16 to 2,048 as the number of nodes in this single layer. We did not observe significant differences between these architectures, so we set the number of nodes to be 1,024 to keep the number close to the initial number of nodes in the imported CNN architectures, as these were initially used to perform classification between 1,000 categories.

We tested two different activation functions for the activation functions of the fully connected layers we added in the side neural network and before the final linear layer. We did not observe any significant differences between the rectified linear units [ReLU] ⁶⁰ and the scaled exponential linear units [SELU] ⁶¹ as activation functions, so we used the more common ReLU.

We then tested different levels of data augmentation. We introduced a hyperparameter that we called “data augmentation factor”. The data augmentation factor modulates the amount of variation introduced by the data augmentation, while preserving the ratio between the different transformations. For example, a data augmentation factor of one is equivalent to the default data augmentation (see Preprocessing - Data augmentation - Images), but a data augmentation factor of two will double the ranges of the possible values sampled and the expected values for the vertical shift, the horizontal shift, the rotation and the zoom on the original images. We tested the following values for the data augmentation factor: 0, 0.1, 0.5, 1, 1.5 and 2. We found that different values for the data augmentation factor hyperparameter yielded similar results, as long as the data augmentation factor was not zero. We therefore set the data augmentation factor to be one when training the final models.

We then tuned the dropout rate for the fully connected layers we added. We tested the following values: 0, 0.1, 0.25, 0.3, 0.5, 0.75, 0.9 and 0.95. We observed that a dropout rate of 0.95 led to underfitting and that smaller values reduced overfitting on the training set but without improving the validation performance. As a consequence, we used a dropout rate of 0.5.

Finally, we tuned the weight decay. We tested the following values: 0, 0.1, 0.2, 0.3, 0.4, 0.5, 1, 5, 10 and 100. For the larger datasets, we found that weight decay values as low as 0.4 could lead to underfitting. We found that lower weight decay values reduced overfitting on the training set without significantly improving the validation performance. We set the weight decay to 0.1.

Altogether, we found that hyperparameter tuning had little effect on the validation performance as long as extreme hyperparameters values were not selected.

Cross-validation

Training deep convolutional neural networks on images and videos is too time and resource consuming to perform a nested cross-validation. Therefore, we tuned the hyperparameters during the preliminary analysis, as described above. After hyperparameters tuning, we performed a simple outer cross-validation to obtain a testing prediction for each sample of the datasets, but we replaced the inner cross-validation with a simple split between the training fold and the validation fold (Table S28). Although the hyperparameters were already tuned, a validation set was still required for two reasons: (1) to perform early stopping⁶², a form of regularization. (2) to generate a set of validation predictions that are necessary for efficient ensemble building (see Methods - Models ensembling) and model selection. During the cross-validation, we scaled and centered the target variable (chronological age) as well as the side predictors (sex and ethnicity) around zero with a standard deviation of one, using the training summary statistics. Scaling the target and the input helps prevent the issues of exploding and vanishing gradients^{63,64}.

Cross-validation example

For the sake of clarity, let us walk through an example. Let us say that we want to generate unbiased predictions for every sample in a dataset using a CNN. First, we select the data fold #0 as the validation set, the data fold #1 as the testing set, and the remaining data folds (#2-9) as the training set. Then we scale and center the target (age), and the side predictors (sex and ethnicity) using the training mean and standard deviation: for each of the variables, we subtract the training mean to the variable on both the training, the validation and the testing set, and we divide it by the training standard deviation. We then train the model on the training set until convergence and select the architecture's parameters (also known as "weights") associated with

the epoch that yielded the lowest validation RMSE. We then use the optimal weights to generate validation predictions for the data fold #0 and testing predictions on the data fold #1. Finally, we unscale the validation and testing predictions by multiplying them by the initial age training standard deviation before adding the initial age training mean to them. This completes the first cross-validation fold.

We then reiterate the process, this time using the data fold #1 as the validation set, the data fold #2 as the testing set, and the remaining data folds (#0 and #3-9) as the training set. We use the optimized weights to generate the validation predictions on the data fold #2, and the testing predictions on the data fold #3. We unscale the validation and testing predictions. This completes the second cross-validation fold. We reiterate the process eight more times to complete the cross-validation. We then concatenate the validation predictions from the ten data folds to obtain the final validation predictions, and the testing predictions from the ten data folds to obtain the final testing predictions.

Interpretability of the predictions

Scalar data-based predictors

For elastic nets, we interpreted the models using the values of the regression coefficients. Large absolute values for these coefficients means they played an important role when generating the predictions. For gradient boosted machines we used the feature importances, which are based on the number of times a tree selected each of the variables. Variables with high feature importances were selected more often and are therefore likely to play a key role in predicting chronological age. For neural networks, we estimated the importance of each feature by

permuting it randomly between samples before computing the performance of the model. The score of each feature is the difference between the R-Squared value before and after the random permutations. Features whose random permutation leads to a large decrease in the model's performance are estimated to be important predictors of chronological age.

We estimated the concordance between the three different algorithms by computing the Pearson and the Spearman correlations between their feature importances.

Image-based predictors

To interpret the CNNs built on images, we first used saliency maps⁶⁵, which we coded using the keract python library. For each input sample, a saliency map uses the gradient of the final prediction with respect to each individual input pixel to estimate whether changing the value of this pixel would affect the prediction. Pixels for which the gradient is close to zero are not important, whereas pixels with a large gradient are estimated to be important.

We then built a second attention map using a custom version of the Gradient-weighted Class Activation Mapping [Grad-CAM] algorithm⁶⁶ adapted to regression rather than multi-class classification: Gradient-weighted Regression Activation Mapping [Grad-RAM]. The intuition behind Grad-CAM maps is that they are similar to saliency maps⁶⁶, but instead of computing the gradient with respect to the input image, they compute it with respect to the activation of the last convolutional layer. As convolutional layers maintain the spatial organization of the input image, Grad-CAM can still identify which region of the image is driving the predictions. Because Grad-CAM does not have to backpropagate the gradient all the way back to the input image, it is considered a less noisy alternative to the saliency maps. In the same way that saliency maps

need to combine the attention maps generated in the different input channels (e.g. RGB) into a single activation map, Grad-CAM must combine the attention maps generated on the different filters of the last convolutional layer. For example, the last convolutional layer for InceptionResNetV2 has 1,792 filters. Grad-CAM combines these 1,792 attention maps into a single attention map using a linear combination. In the initial Class Activation Mapping [CAM] algorithm ⁶⁷, generating CAM activation maps required to retrain the model after modifying the architecture and replacing all the fully connected layers after the final convolutional layer with a global max pooling operation, which converted each filter into a scalar feature. The intuition behind this substitution was that each filter could be interpreted as detecting a specific feature, and global max pooling yielded a scalar that could be interpreted as the presence (high value) or absence (low value) of the feature anywhere on the image. The scalar values were then linearly combined and activated using the softmax function to yield the probabilities of belonging to different classes. To obtain the activation map for a specific class, the filters of the last convolution layer were linearly combined using the weights connecting the scalar features obtained after the max pooling operation to the final prediction score for that class. CAM was later improved to become Grad-CAM ⁶⁶. Grad-CAM saves the need for modifying the architecture of the model and retraining it by approximating the linear regression weight for each final convolutional filter by the mean activation gradient over the pixels of the filter. The intuition behind this approximation is that a filter's pixel is important if changing its value affects the final prediction, so a high average gradient over the pixels of the filter justifies that this filter should be given a higher weight when merging all the filters into a single attention map. To adapt Grad-CAM to our regression task we (1) computed the derivatives of the chronological age prediction rather than a class' prediction and (2) removed the ReLU activation applied to the weighted sum of the last convolutional filters, which we replaced by an absolute value. The

rationale is that for (Grad-)CAM maps, we only want to highlight the regions of the picture which are associated with a high probability for the class. In contrast, for (Grad-)RAM we care as much about the regions of the input image that can strongly increase the chronological age prediction as about the regions that can strongly decrease it. Because the filters in the last convolutional layer are the result of the processing of the input image by several convolutional layers with possibly negative weights, the sign of the last convolutional layer's pixels and regression weights cannot be linked to either accelerated aging or decelerated aging, only to the magnitude of the shift that would affect the prediction if each region of the input image was modified. Regression Activation Mapping (RAM) was mentioned as a possible extension of CAM in the original CAM publication ⁶⁷ and has been used to interpret models CNNs built on retinal images ⁶⁸ and cortical surfaces ⁶⁹, but we are to our knowledge the first to describe the generalization of Grad-CAM to a regression task. One notable difference between our implementation and Wang and Yang.'s implementation ⁶⁸ is that we are taking the absolute value of the final attention map, as mentioned above. We found that not taking the absolute value led to misleading attention maps for participants with high chronological age predictions. The attention map highlights important areas with negative values, which are therefore depicted in blue, a color otherwise associated with unimportant regions in traditional CAMs. Inversely, regions on the input image for which the attention map has a slight positive value are spuriously considered to be the most important and are highlighted in red. We therefore advise that RAM or Grad-RAM be implemented using an absolute value. We coded Grad-RAM using the `get_activations` and `get_gradients_of_activations` functions of the `keract` python library.

It is important to understand that unlike the feature importances described under "Scalar data-based predictors", which describe the model itself, attention maps are sample specific. In

other words, they can be used to explain which features drove the predictions for a specific inputted sample but cannot provide an explanation for the way the model is performing predictions in general.

For each aging subdimension, we generated the attention maps for the best performing CNN architecture. We selected representative samples for which we computed the different attention maps. We computed attention maps for the two sexes (female and male), for three age ranges (ten youngest ages, ten middle ages and ten oldest ages of the chronological age distribution) and for three aging rates (accelerated agers, normal agers, decelerated agers). For each intersection of the three categories listed above, we selected the ten most representative samples (e.g. the ten most accelerated agers among young males). The figures in this paper only present the first, most representative of these ten samples. The complete set of samples can be found on the website.

Non-genetic correlates of accelerated aging

Unlike DNA, biomarkers, phenotypes, diseases, family history, environmental variables and socioeconomics can change over life. As a consequence, we compared each biomarker, phenotype and environmental variable with the accelerated aging of the participant at the time the exposure was measured and we used the “Samples predictions”, as opposed to the “Participants predictions” that we used for the identification of genetic correlates (see Methods - Models ensembling - Generating average predictions for each participant).

Imputation of the non-genetic X-variables

Most X-variables were not collected on all four instances. Additionally, no X-variables were collected at the same time as the accelerometer data was collected. To identify the non-genetic correlates of accelerated aging, we had to impute the values of the X-variables for the ages of the participants for which they were not available. We considered two imputation methods, which we refer to as the “cross-sectional” and the “longitudinal” imputations.

For the cross-sectional imputation, we computed a linear regression for each X variable as a function of age, adjusting for sex. We then used the slope of the linear regression to extrapolate the value of the XWAS variable at different ages.

For the longitudinal imputation, we first selected, for each X variable, all the participants that had at least two measures taken for this X variable. We then performed a linear regression for each participant. We then averaged the slope of the linear regressions over all the participants of the same sex. Finally, we used this slope to extrapolate the value of the XWAS variable at different ages for all participants depending on their sex, in the same way we did it for the cross-sectional imputation.

It is important to notice that for both the cross-sectional imputation and the longitudinal imputation, data can only be imputed when the XWAS variable has been measured at least once for the participant. This raw measure is then used to extrapolate which value the X variable was likely taking a couple years earlier and/or later.

The advantage of the cross-sectional imputation is larger sample sizes. The advantage of the longitudinal method is that it corrects for generational effects. For example, old people have shorter legs than young people on average ⁷⁰. This is not because human legs shrink as we grow older. Instead, people who are old today already had shorter legs when they were young. If the cross-sectional regression is used to impute the length of the participants on instances where it was not measured, it will spuriously assign smaller values to the older samples. In contrast, the longitudinal regression learns the regression coefficient by comparing each participant to themselves as they age and will therefore not capture the generational effect. When used to predict the participants legs' length, it will impute constant values over time. To evaluate which of the two imputation methods should be preferred, we used them to predict X-variables for which we knew the actual values and computed the R-Squared values associated with the predictions. We found that, even with sample sizes as small as 200 samples, longitudinal imputation outperformed cross-sectional imputation. We therefore used longitudinal imputation.

X-Wide Association Studies

First, we tested for associations in an univariate context by computing the partial correlation between each X-variable and musculoskeletal aging dimensions. To compute the partial correlation between an X-variable and an aging, we followed a three steps process. (1) We ran a linear regression on each of the two variables, using age, sex and ethnicity as predictors. (2) We computed the residuals for the two variables. (3) We computed the correlation between the two residuals and the associated p-value if their intersection had a sample size of at least ten samples. We used a threshold for significance of 0.05 and corrected the p-values for multiple

testing using the Bonferroni correction. We plotted the results using a volcano plot. We refer to this pipeline as an X-Wide Association study [XWAS].

Supplementary Figures

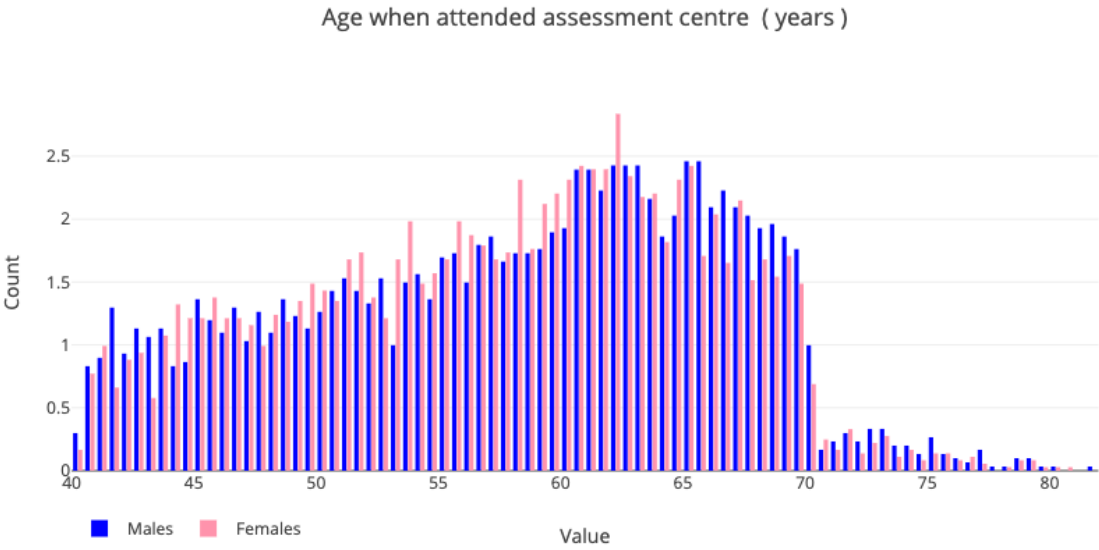


Figure S1: Demographics of the UK Biobank cohort

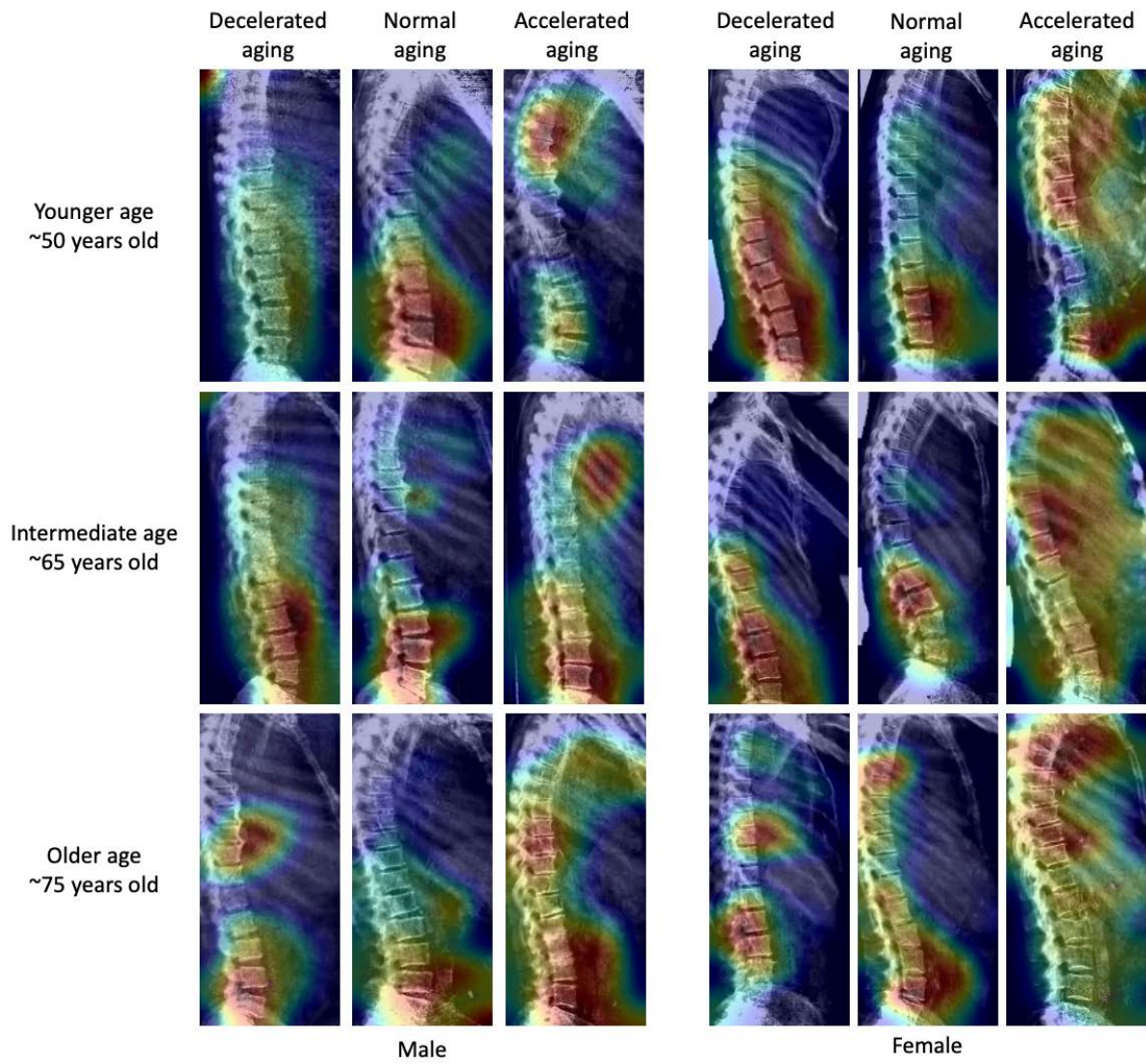


Figure S2: Attention maps for spine X-ray images (sagittal view)

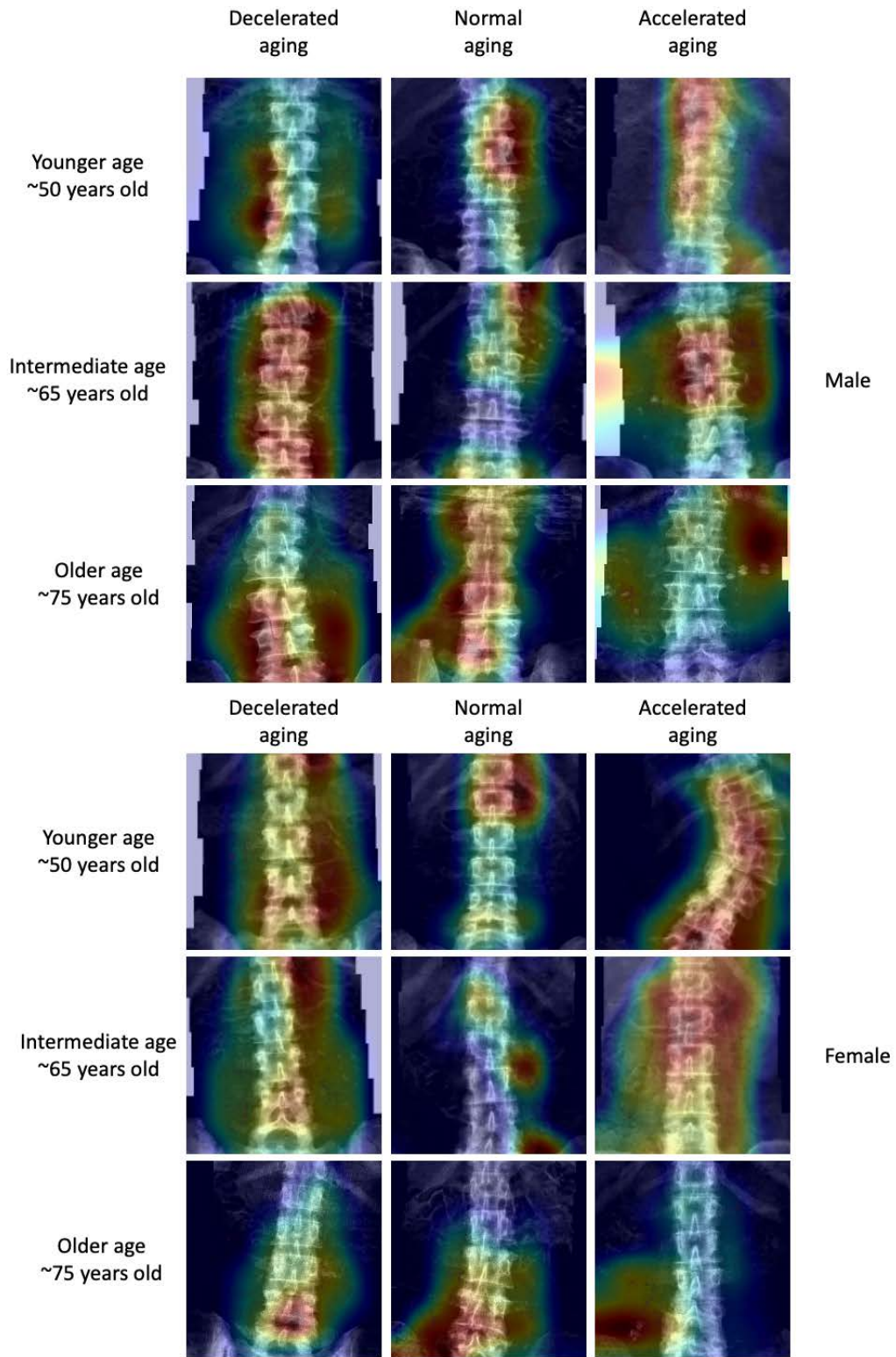


Figure S3: Attention maps for spine X-ray images (coronal view)

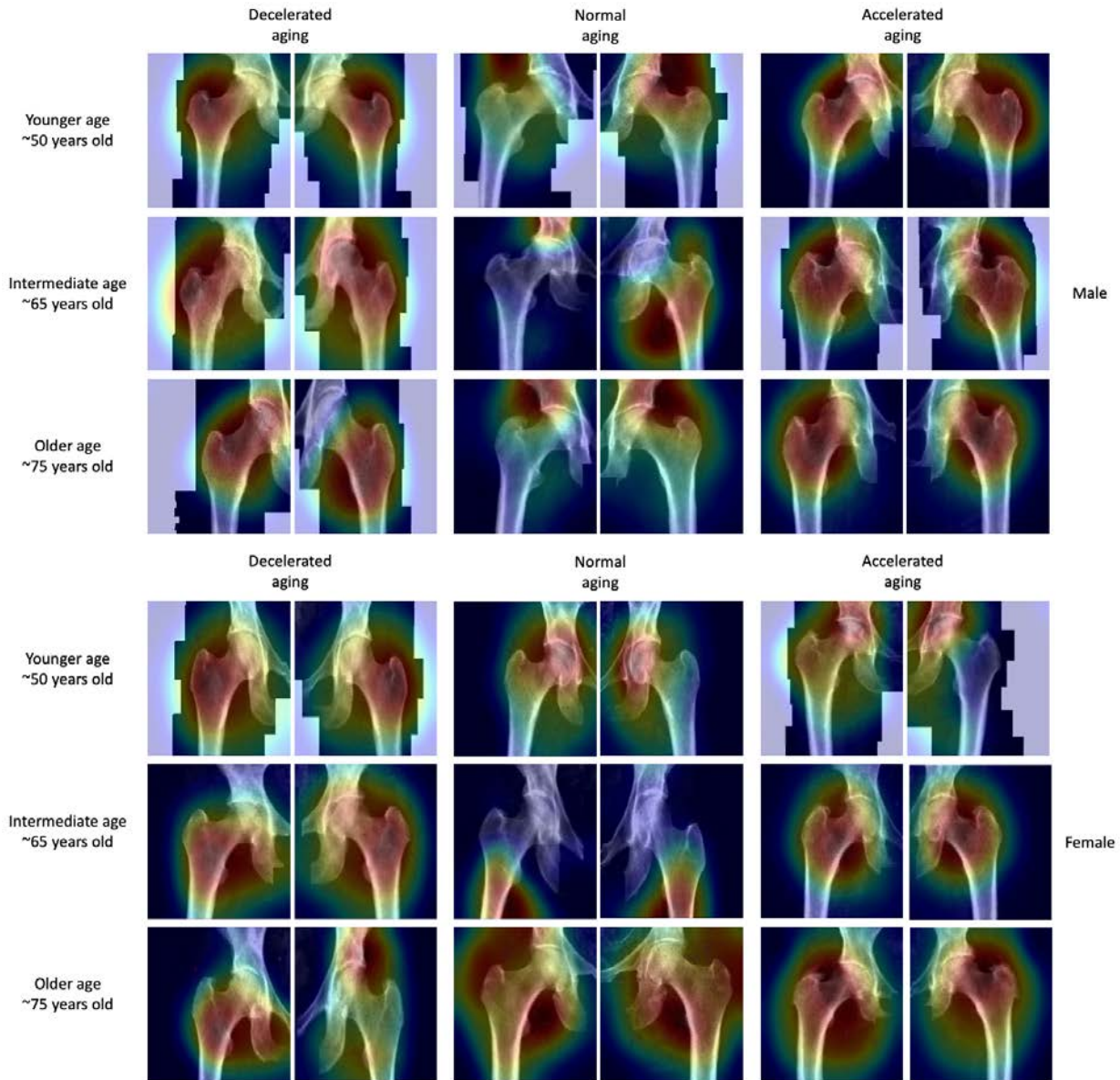


Figure S4: Attention maps for hip X-ray images

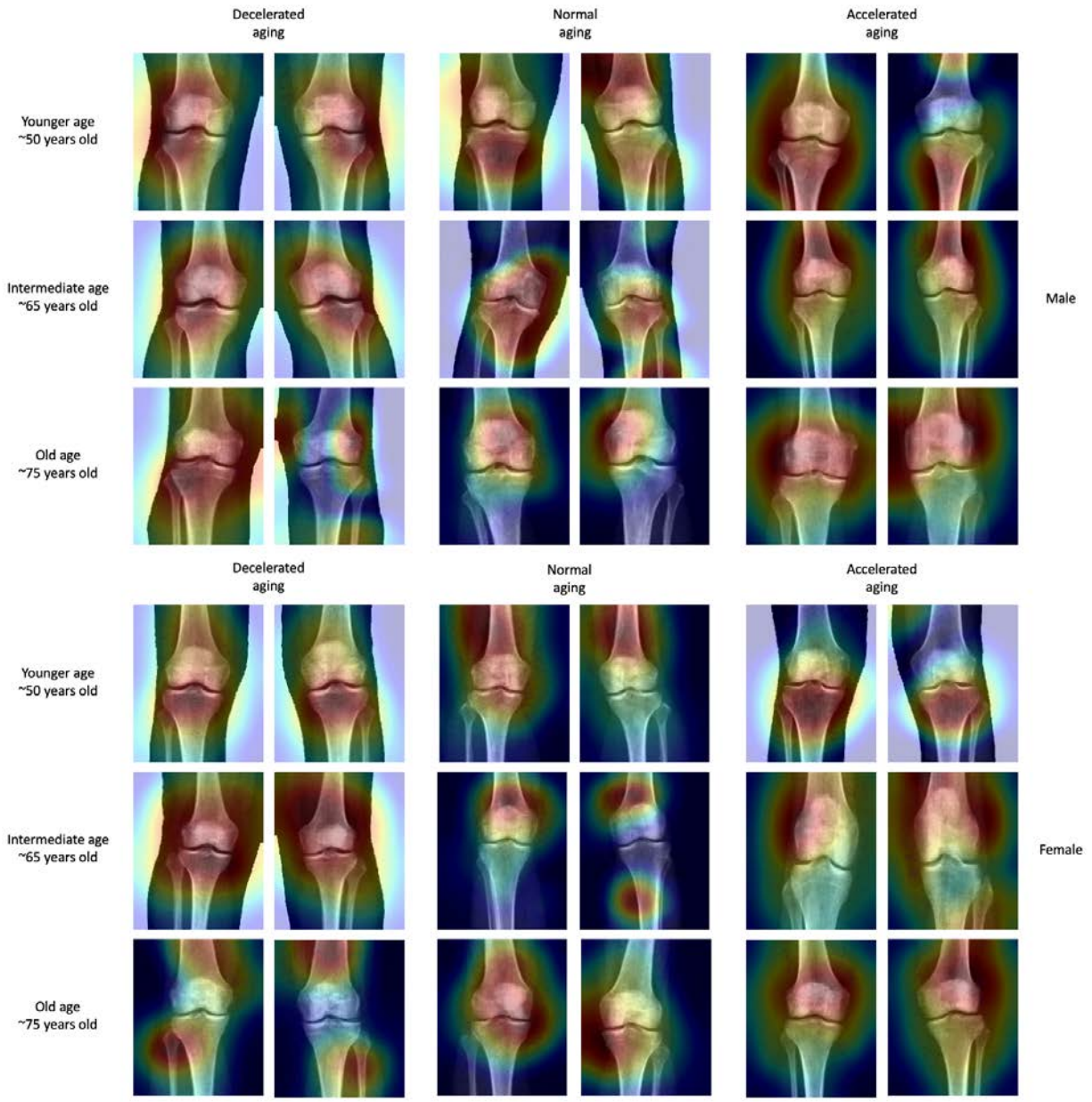


Figure S5: Attention maps for knee X-ray images

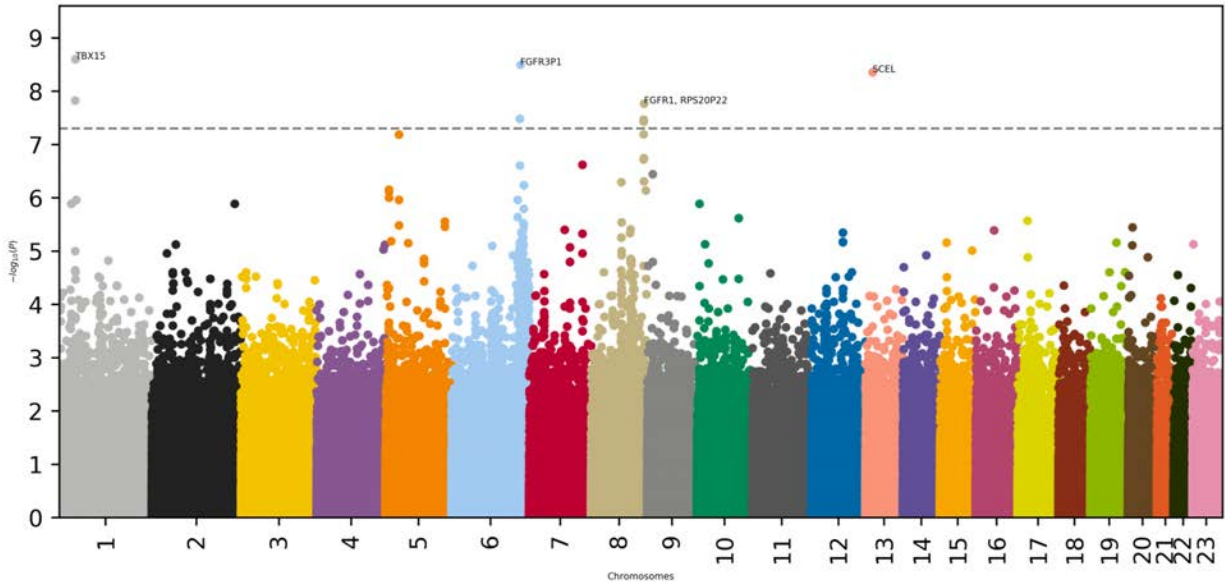


Figure S6: GWAS results - Full body aging

$-\log_{10}(p\text{-value})$ vs. chromosomal position of locus. Dotted line denotes 5×10^{-8} .

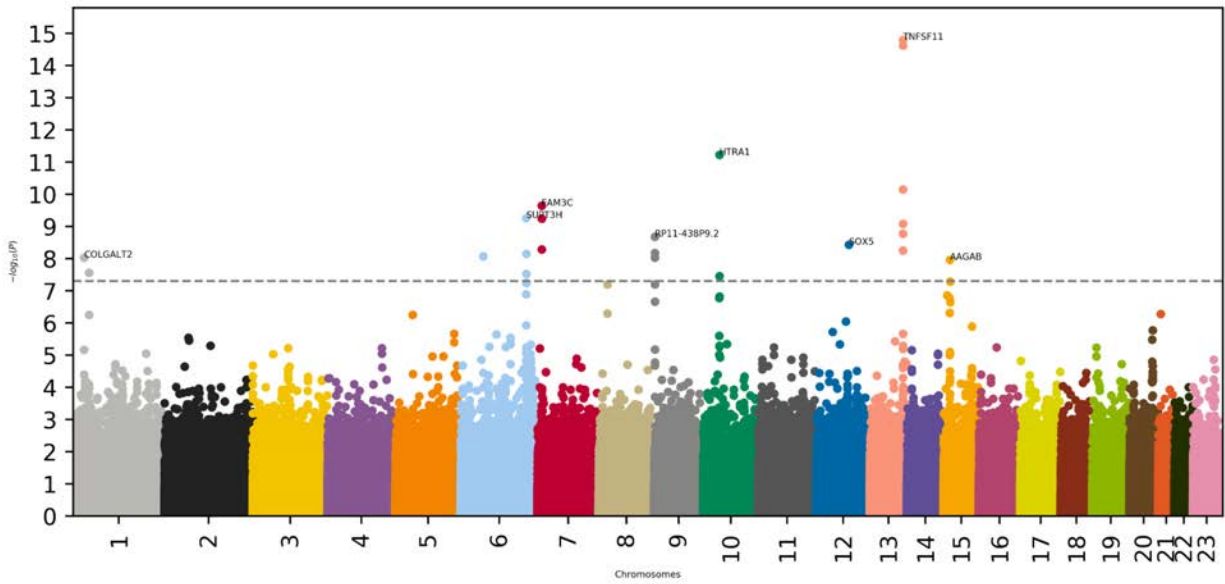


Figure S7: GWAS results - Spine aging

$-\log_{10}(p\text{-value})$ vs. chromosomal position of locus. Dotted line denotes 5×10^{-8} .

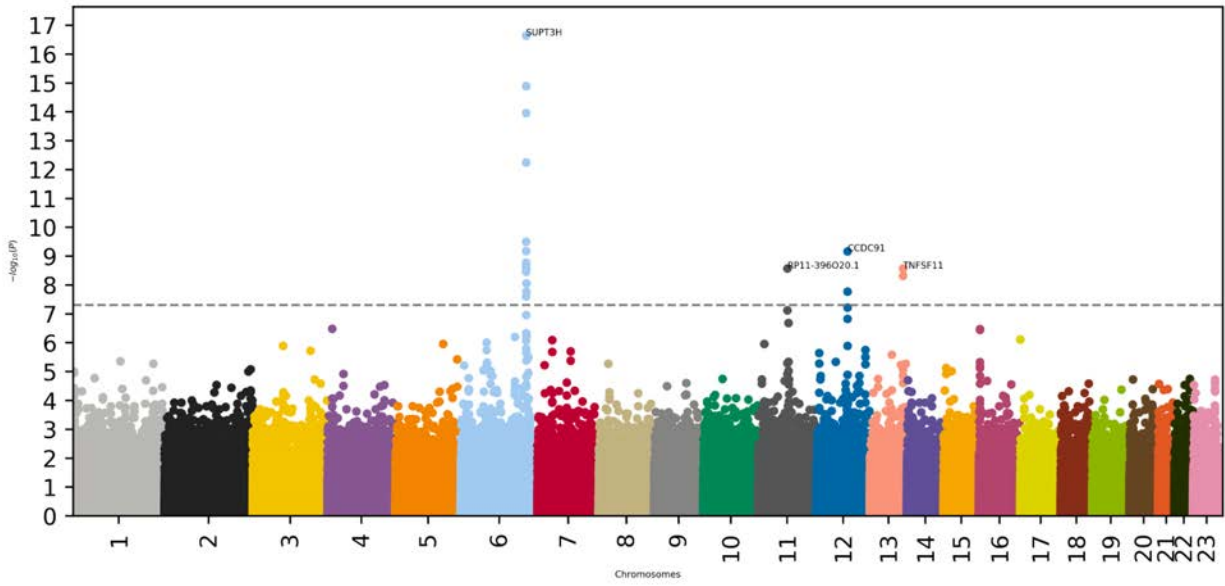


Figure S8: GWAS results - Hip aging

$-\log_{10}(p\text{-value})$ vs. chromosomal position of locus. Dotted line denotes 5×10^{-8} .

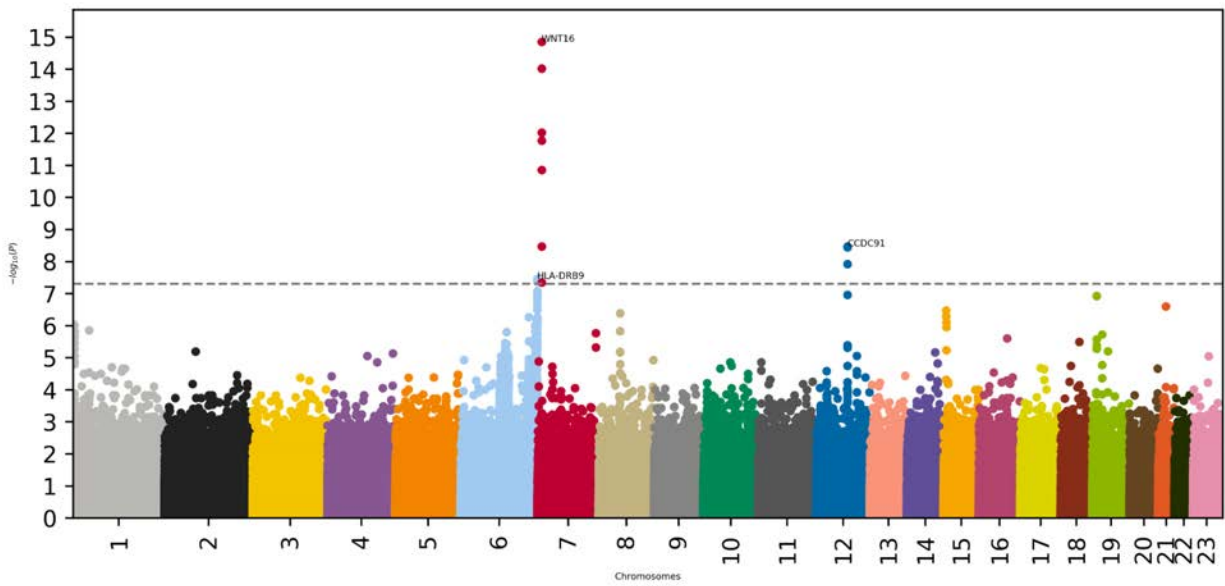


Figure S9: GWAS results - Knee aging

$-\log_{10}(p\text{-value})$ vs. chromosomal position of locus. Dotted line denotes 5×10^{-8} .

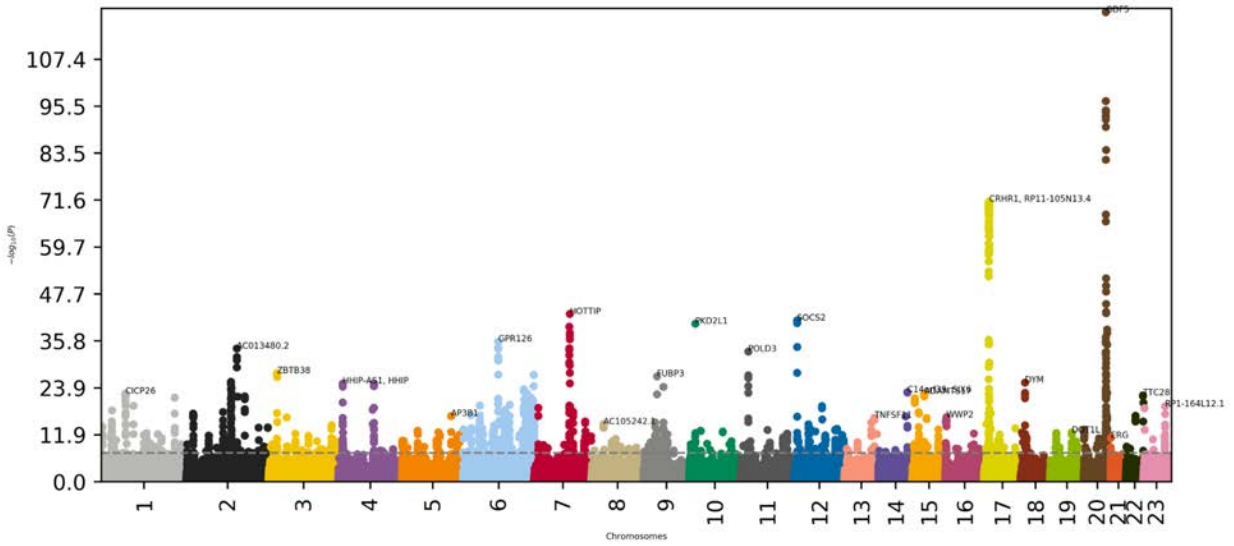


Figure S10: GWAS results - Anthropometry, impedance, heel bone densitometry and hand grip strength-based aging

$-\log_{10}(p\text{-value})$ vs. chromosomal position of locus. Dotted line denotes 5×10^{-8} .

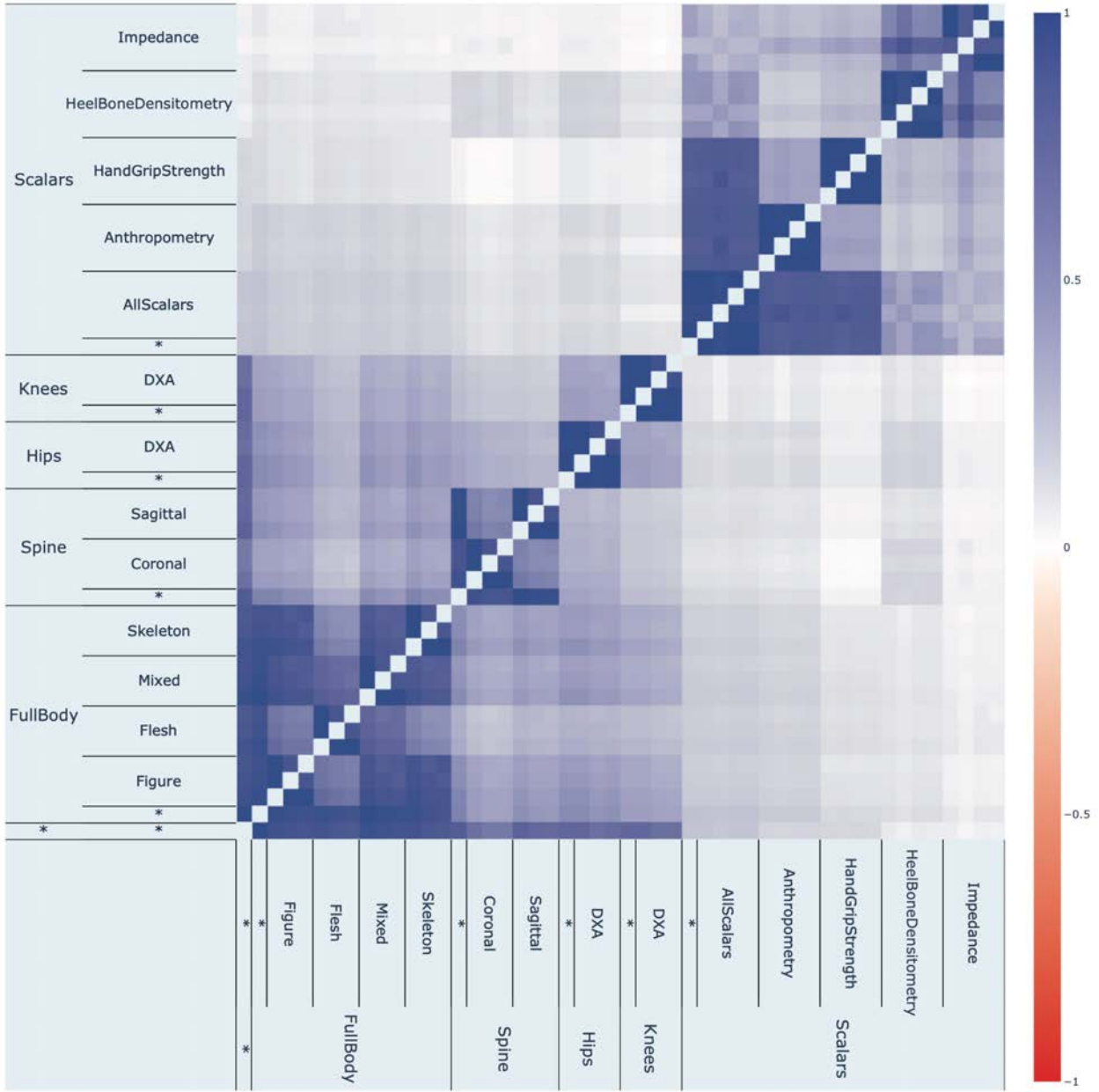


Figure S11: Phenotypic correlation between all accelerated musculoskeletal aging dimensions and subdimensions

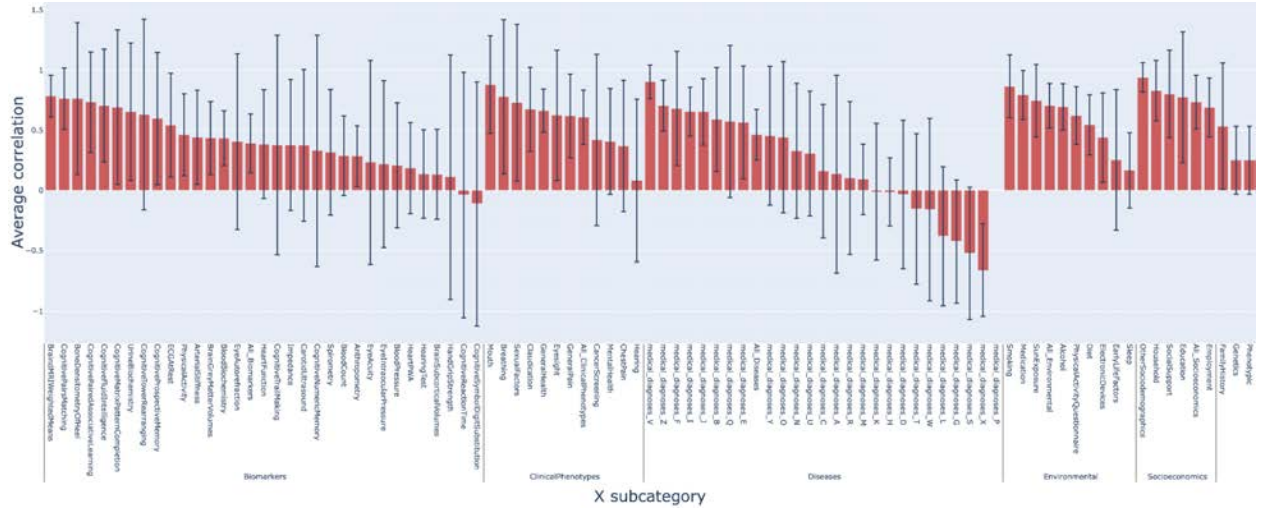


Figure S12: Average correlation between musculoskeletal accelerated aging dimensions in terms of association with non-genetic factors

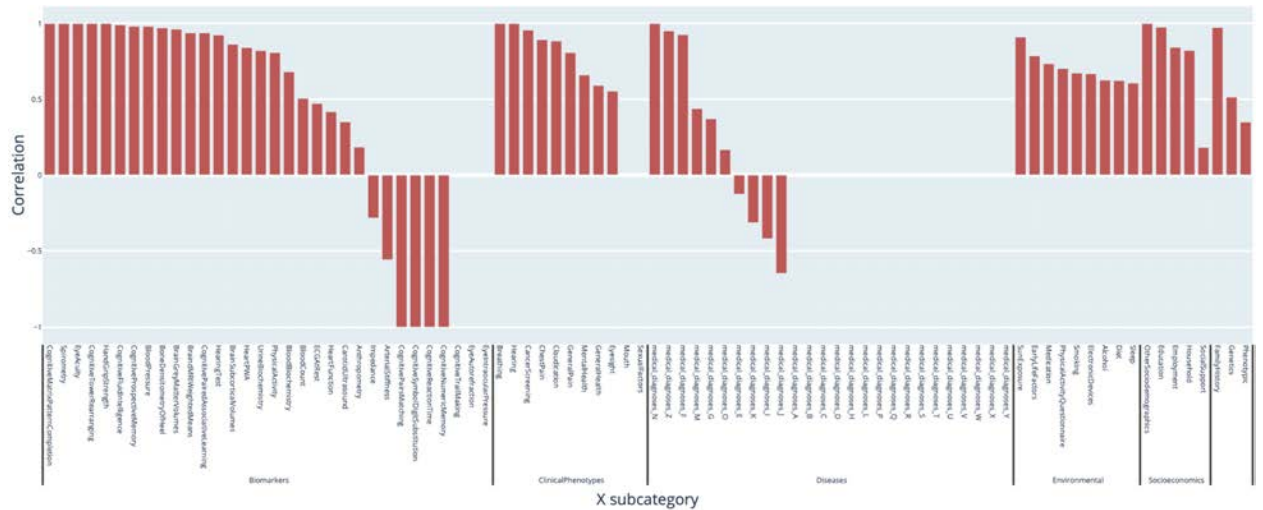


Figure S13: Correlation between accelerated hip and knee aging in terms of associated biomarkers, clinical phenotypes, diseases, family history, environmental and socioeconomic variables

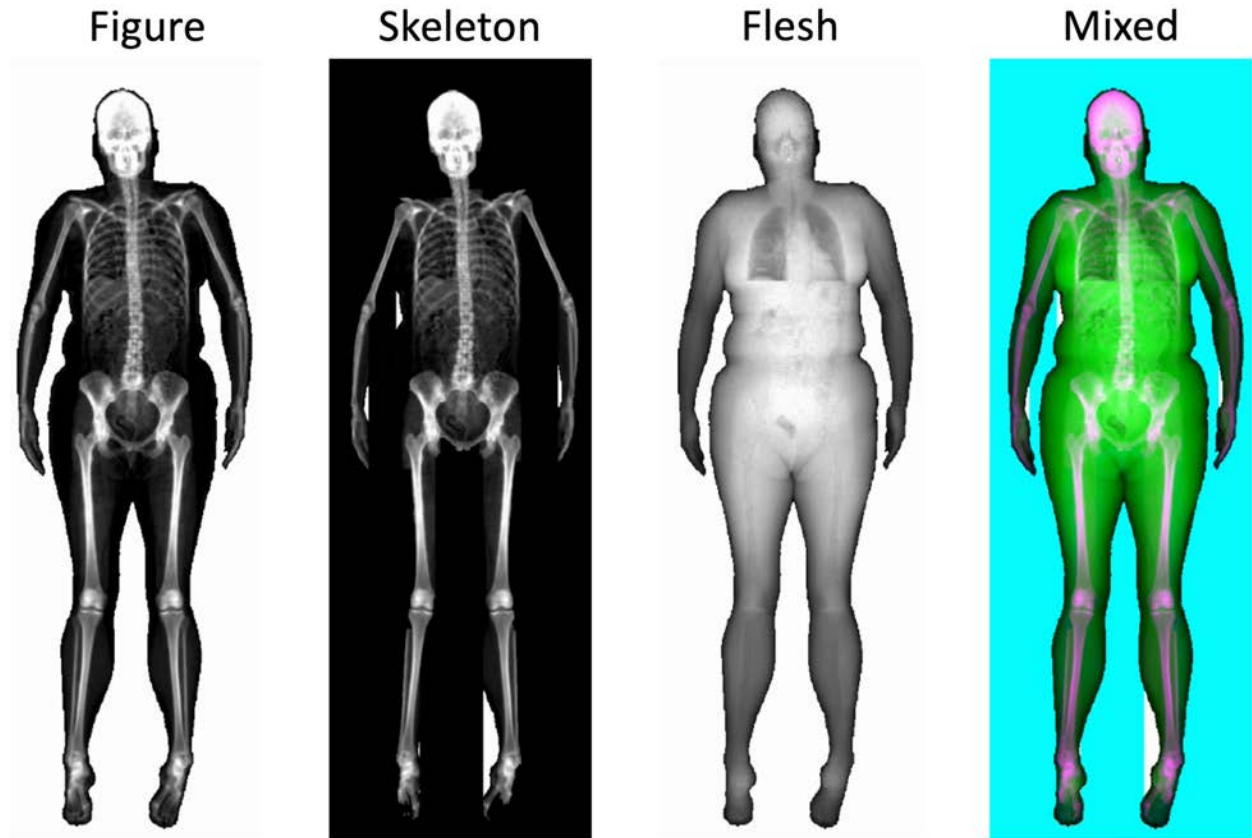


Figure S14: Sample preprocessed musculoskeletal full body X-ray images

The participant is a 50-55-year-old female.

Sagittal



Coronal

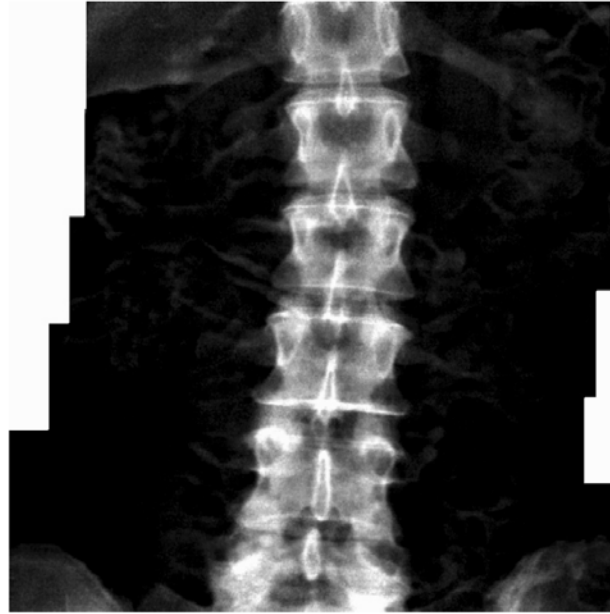


Figure S15: Sample preprocessed spine X-ray images

The participant is a 50-55-year-old female.

Left



Right

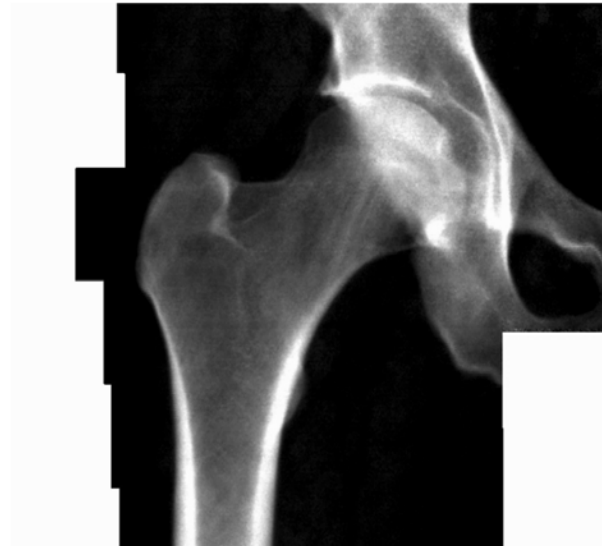


Figure S16: Sample preprocessed hip X-ray images

The participant is a 60-65-year-old male.

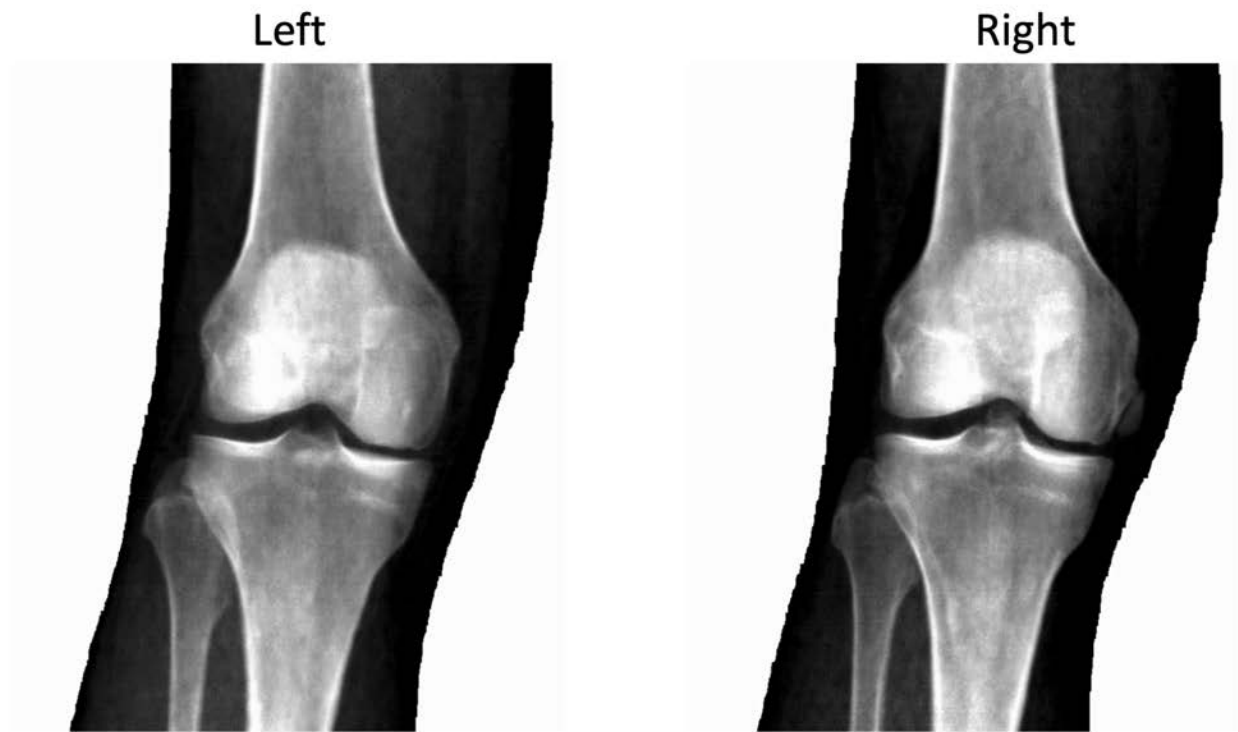


Figure S17: Sample preprocessed knee X-ray images

The participant is a 60-65-year-old male.

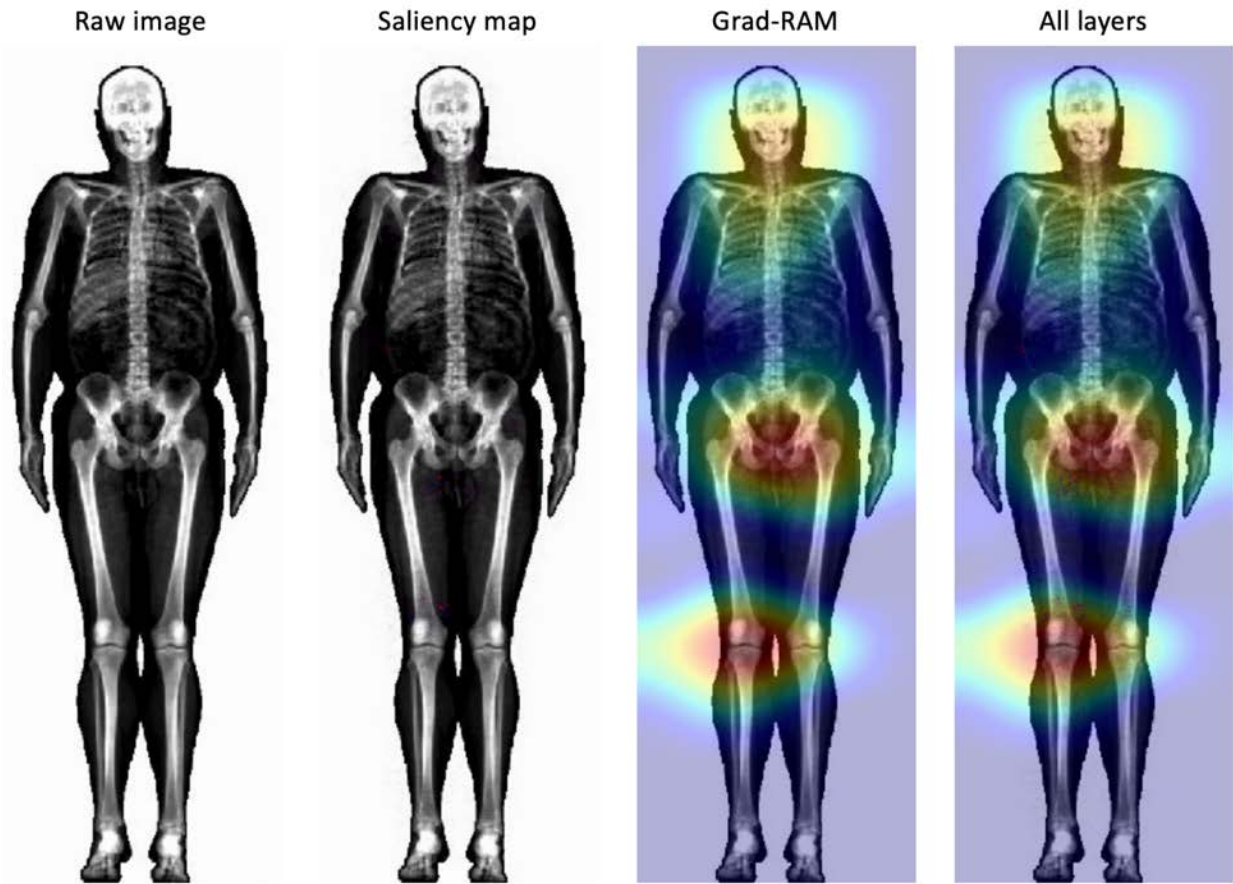


Figure S18: Attention map algorithms: saliency map and Grad-RAM.

The participant is a 60-65-year-old male whose chronological age was accurately predicted by the machine using a “Figure” musculoskeletal X-ray image.

Supplementary Tables

Table S1: Comparison between our musculoskeletal age predictors and the literature in terms of prediction performance

Musculoskeletal aging dimension	Our model				Model(s) in the literature					
	R-Squared (%)	RMSE (years)	Sample size	Age range (years)	R-Squared (%)	RMSE (years)	Algorithm	Sample size	Age range (years)	Authors
Spine	74.6±0.2	3.81±0.01	40,951	45.1-82.5	N/A					
Hips	69.0±0.3	4.20±0.02	83,234	45.1-82.5	38	1.91	Linear regression	643	10-30	Wittschieber et al.
Knees	69.0±0.3	4.20±0.02	79,477	45.1-82.5	77.5-81.5; 88-90; 44.1-65.4	N/A; N/A; N/A	Linear regression; Linear regression; Linear regression	221; 503; 322	9-19; 6-19; 11-32	O'Connor et al.; Tang et al.; Fan et al.
Full body	85.7±0.1	2.85±0.01	42,164	45.1-82.5	82	2.47±1.91 (MAE)	CNN (VGG16)	32,323	44-82	Langner et al.
Hand	N/A				N/A	1.01 (MAE)	CNN (Xception)	12,811	0-19	Westerberg
Scalar-Biomarkers	25.9±0.01	7.10±0.01	486,642	37.4-82.3	N/A					

Table S2: Comparison between the models trained on scalar features

Musculoskeletal subdimension	Number of Predictors (non-demographics)	ElasticNet (R-Squared)	LightGBM (R-Squared)	NeuralNetwork (R-Squared)
Impedance	5	0.020±0.001	0.042±0.002	0.041±0.005
Anthropometry	8	0.132±0.003	0.151±0.003	0.147±0.004
HeelBoneDensitometry	6	0.035±0.001	0.049±0.002	0.048±0.002
HandGripStrength	2	0.117±0.003	0.131±0.003	0.130±0.003
AllScalars	21	0.205±0.004	0.252±0.004	0.252±0.005

Table S3: Pearson correlations between the feature importances for different scalar features-based algorithms

Musculoskeletal subdimension	Correlation Versus ElasticNet	Correlation Versus LightGBM	Correlation Versus NeuralNetwork	ElasticNet Versus LightGBM	ElasticNet Versus NeuralNetwork	LightGBM Versus NeuralNetwork
------------------------------	-------------------------------	-----------------------------	----------------------------------	----------------------------	---------------------------------	-------------------------------

Impedance	0.535	0.043	0.72	0.385	0.344	0.032
Anthropometry	0.402	0.361	0.389	0.712	0.648	0.288
HeelBoneDensitometry	0.611	0.699	0.418	0.499	0.446	0.033
HandGripStrength	0.581	0.837	0.49	0.706	0.584	0.277
AllScalars	0.337	0.352	0.241	0.588	0.495	0.213

Table S4: Spearman correlations between the feature importances for different scalar features-based algorithms

Musculoskeletal subdimension	Correlation Versus ElasticNet	Correlation Versus LightGBM	Correlation Versus NeuralNetwork	ElasticNet Versus LightGBM	ElasticNet Versus NeuralNetwork	LightGBM Versus NeuralNetwork
Impedance	0.64	0.49	0.796	0.582	0.455	0.608
Anthropometry	0.329	0.563	0.588	0.685	0.523	0.739
HeelBoneDensitometry	0.526	0.732	0.457	0.638	0.232	0.357
HandGripStrength	0.461	0.693	0.545	0.557	0.477	0.703
AllScalars	0.331	0.513	0.384	0.566	0.567	0.618

Table S5: List of biomarkers by subcategories for the Biomarkers Wide Association Study [BWAS]

See supplementary data

Table S6: Biomarkers most associated with accelerated aging for each musculoskeletal aging dimension

See supplementary data

Table S7: Biomarkers most associated with decelerated aging for each musculoskeletal aging dimension

See supplementary data

Table S8: List of clinical phenotypes by subcategories for the Clinical Phenotypes Wide Association Study [CWAS]

See supplementary data

Table S9: Clinical phenotypes most associated with accelerated aging for each musculoskeletal aging dimension

See supplementary data

Table S10: Clinical phenotypes most associated with decelerated aging for each musculoskeletal aging dimension

See supplementary data

Table S11: List of diseases by subcategories for the Diseases Wide Association Study [DWAS]

See supplementary data

Table S12: Diseases most associated with accelerated aging for each musculoskeletal aging dimension

See supplementary data

Table S13: Diseases most associated with decelerated aging for each musculoskeletal aging dimension

See supplementary data

Table S14: List of family history variables by subcategories for the Family History Phenotypes Wide Association Study [FWAS]

See supplementary data

Table S15: Family history variables most associated with accelerated aging for each musculoskeletal dimension

Musculoskeletal dimension	X-subcategory	Variables
General	FamilyHistory (0.0%)	
Full body	FamilyHistory (1.1%)	
Hip	FamilyHistory (1.1%)	
Knee	FamilyHistory (0.0%)	
Scalar biomarkers	FamilyHistory (3.3%)	Illnesses of siblings.Prefer not to answer1 (.013); Illnesses of siblings.Severe depression (.011); Illnesses of mother.Severe depression (.009)
Spine	FamilyHistory (0.0%)	

Table S16: Family history variables most associated with decelerated aging for each musculoskeletal dimension

Musculoskeletal dimension	X-subcategory	Variables
General	FamilyHistory (1.1%)	
Full body	FamilyHistory (0.0%)	
Hip	FamilyHistory (1.1%)	Illnesses of father.Prostate cancer (.027)
Knee	FamilyHistory (3.3%)	Illnesses of father.Prostate cancer (.025); Illnesses of mother.Severe depression (.025)
Scalar biomarkers	FamilyHistory (6.5%)	Illnesses of father.Lung cancer (.141); Illnesses of adopted mother.Chronic bronchitis/emphysema (.124); Illnesses of adopted father.Lung cancer (.121)
Spine	FamilyHistory (1.1%)	

Table S17: List of environmental variables by subcategories for the Environmental Wide Association Study [EWAS]

See supplementary data

Table S18: Environmental variables most associated with accelerated aging for musculoskeletal each aging dimension

See supplementary data

Table S19: Environmental variables most associated with decelerated aging for each musculoskeletal aging dimension

See supplementary data

Table S20: List of socioeconomic variables by subcategories for the Socioeconomics Wide Association Study [SWAS]

See supplementary data

Table S21: Socioeconomic variables most associated with accelerated aging for each musculoskeletal aging dimension

See supplementary data

Table S22: Socioeconomic variables most associated with decelerated aging for each musculoskeletal aging dimension

See supplementary data

Table S23: Image sizes after resizing

Musculoskeletal Dimension	Musculoskeletal Subdimension	Size before resizing	Size after resizing
Spine	Sagittal	1513, 684	466, 211
	Coronal	724, 720	315, 313

Hips	MRI	626, 680	329, 303
Knees	MRI	851, 700	347, 286
FullBody	Mixed	811, 272	541, 181
	Figure	811, 272	541, 181
	Skeleton	811, 272	541, 181
	Flesh	811, 272	541, 181

Table S24: Data augmentation hyperparameters for X-ray images

Dataset	Rotation range	Horizontal shift percentage range	Vertical shift percentage range	Zoom range
Musculoskeletal Spine	10	0.1	0.1	0
Musculoskeletal Hips	10	0.1	0.1	0.1
Musculoskeletal Knees	10	0.1	0.1	0.1
Musculoskeletal FullBody	10	0.05	0.02	0

Table S25: Hyperparameter space for scalar features-based models Bayesian optimization

Algorithm	Hyperparameter	Scale	Low	High
Elastic net	alpha	loguniform	-10	0
	l1_ratio	uniform	0	1
Gradient Boosted Machine	num_leaves	quniform	5	45
	min_child_samples	quniform	100	500
	min_child_weight	loguniform	-5	4
	subsample	uniform	0.2	0.8
	colsample	uniform	0.4	0.6
	reg_alpha	loguniform	-2	2
	reg_lambda	loguniform	-2	2
n_estimators	quniform	150	450	
Neural network	learning_rate_init	loguniform	-5	-1
	apha	loguniform	-6	3

Table S26: Nested Cross-Validation pipeline

10x9-fold Nested Cross-Validation [CV] 10-fold Outer Cross-Validation; 9-fold Inner Cross-Validation (N_CV_inner = N_CV_outer - 1)												
INITIAL DATA SPLIT	Data Fold	N=100										
	F0	N=10										
	F1	N=10										
	F2	N=10										
	F3	N=10										
	F4	N=10										
	F5	N=10										
	F6	N=10										
	F7	N=10										
	F8	N=10										
F9	N=10											
OUTER CROSS-VALIDATION FOLD 0	Data Fold	Outer CV folds 0	Inner CV folds 0	Inner CV folds 1	Inner CV folds 2	Inner CV folds 3	Inner CV folds 4	Inner CV folds 5	Inner CV folds 6	Inner CV folds 7	Inner CV folds 8	
	F0	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F1	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F2	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F3	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F4	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F5	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F6	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F7	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F8	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
F9	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train		
OUTER CROSS-VALIDATION FOLD 1	Data Fold	Outer CV folds 1	Inner CV folds 0	Inner CV folds 1	Inner CV folds 2	Inner CV folds 3	Inner CV folds 4	Inner CV folds 5	Inner CV folds 6	Inner CV folds 7	Inner CV folds 8	
	F0	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F1	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F2	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F3	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F4	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F5	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F6	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F7	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F8	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
F9	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train		
OUTER CROSS-VALIDATION FOLD 9	Data Fold	Outer CV folds 9	Inner CV folds 0	Inner CV folds 1	Inner CV folds 2	Inner CV folds 3	Inner CV folds 4	Inner CV folds 5	Inner CV folds 6	Inner CV folds 7	Inner CV folds 8	
	F0	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F1	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F2	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F3	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F4	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F5	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F6	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F7	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
	F8	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	
F9	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train		
SPECIAL "ALL FOLDS" PERMUTATION for feature importance	Data Fold	Outer CV folds All	Inner CV folds 0	Inner CV folds 1	Inner CV folds 2	Inner CV folds 3	Inner CV folds 4	Inner CV folds 5	Inner CV folds 6	Inner CV folds 7	Inner CV folds 8	Inner CV folds 9
	F0	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train
	F1	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train
	F2	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train
	F3	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train
	F4	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train
	F5	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train
	F6	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train
	F7	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train
	F8	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train
F9	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	Train, Train	

METHODS:

General comments: The Nested Cross-Validation is a normal double split train/test method, for which not just one, but the two splits are replaced with a CV instead. The Outer Cross-Validation is used to replace the split between train + val on one side, vs. test on the other side. Thanks to the Outer Cross-Validation, an unbiased testing prediction can be obtained for every sample. Three common sources of confusion clarified: Yes, using an Outer Cross-Validation does mean that the predictions do not all come from the same model. (For example for testing, the prediction for each fold comes from a different model.) No, the Outer Cross-Validation has "feedback" to do with hyperparameter tuning or model selection. Tuning a model only requires one Outer Cross-Validation. In contrast, it requires N_CV_folds distinct Inner Cross-Validations, one for each fold of the Outer Cross-Validation. (see below)

Inner Cross-Validation: The Inner Cross-Validation is used to replace the split between the typical training set and validation set. Thanks to the Inner Cross-Validation, hyperparameter selection can be decided based on a performance calculated on the full training set and validation set sample size, instead of the side validation set sample size. Two common sources of confusion clarified: #1 - The whole purpose of performing an Inner Cross-Validation is "hyperparameter selection". The Inner Cross-Validation tunes the model that will then be used to generate predictions on the left out testing fold. Therefore, the Inner Cross-Validation does "test" generate a testing prediction for more samples. That is the role of the Outer Cross-Validation. #2 - Tuning a model requires N_CV_folds distinct Inner Cross-Validations, one for each fold of the Outer Cross-Validation. In contrast, it only requires one Outer Cross-Validation. (see above)

Special Fold 0: Special "All folds" permutation is only here for feature importance purposes. We could extract the feature importance from the 10 models used to generate the predictions on the testing set, but we would have to take the average of feature importances between several models. Additionally each of these models only saw 90% of the data (10% was left aside for testing purposes). If we only care about feature importance, we actually do not need to leave samples aside for testing. We can use 100% of the data to train a model using a 10-fold (not 9-fold) time, since the test fold can be used too! "Inner" Cross-Validation to train a single model. This model will never be tested, but it does not matter. Its role is to provide the best estimate for the feature importance, which is done well.

Two common sources of confusion clarified: #1 - The Nested Cross-Validation pipeline and the Special "All folds" permutation pipeline are "test" determined and serve different purposes. The purpose of the Nested Cross-Validation pipeline is to generate accurate and unbiased testing prediction for every sample, whereas the purpose of the Special "All folds" pipeline is to generate feature importance. #2 - Yes, it does mean that the feature importance reported do not correspond to any of the 10 models (trained in the 10 Outer Cross-Validation folds) that generated the final testing predictions. The feature importance comes from a model that was never used to generate predictions.

RESULTS:	Example: cells corresponding to the values used for each prediction set	Explanation:
Training	"E1H2I3J4K5L6M7N8O9P10Q11R12S13T14U15V16W17X18Y19Z20"	To obtain the training prediction for the Data Fold F1, take the mean of every training prediction available on this data fold (10*9=90 of them).
Validation	"123"	To obtain the validation prediction for the Data Fold F1, take the prediction outputted by the model trained on the remaining 8 data folds with the best hyperparameters values obtained for this Outer Cross-Validation fold. These validation predictions can then be used to tune the ensemble models. (// Warning //) This can only be properly done if data folds split used for the Outer Cross-Validation and the Inner Cross-Validation is the same. No extra step is needed for the testing prediction for the Data Fold 1. Simply take the only testing prediction available for this fold.
Testing	"234"	
Feature importance:	N/A	Extract the feature importance from the special "outer fold" "All". If want to compute the standard deviation for the feature importance as well, we can leverage the feature importance computed on each of the Outer Cross-Validation folds and take the square root of their variance.

SUPPLEMENTARY:

Comparison between Nested Cross-Validation and other validation methods:

- #1 Single split train/test. Split train/test works when there are no hyperparameters to tune. If one uses a simple train/test split and tune the hyperparameters on the testing set, some overfitting is to be expected.
- #2 Double split train/validation/test. Because of the above, a train/validation/test double split is commonly used when hyperparameters tuning is required. The training is performed on the training set, and the evaluation of the hyperparameters on the validation set, and the evaluation of the unbiased model performance on the untouched test dataset.
- #3 Cross-Validation(train/val) and split (train/val)/test. With the pipeline #2 (double split train/validation/test), the tuning of the hyperparameters is not leveraging the dataset as much as with a Cross-Validation, as only one data fold is used to tune the hyperparameters. Therefore machine learning practitioners usually replace the "inner split" (the training/validation split), by a Cross-Validation instead, but they keep the "outer split" (between the training-validation dataset used for their Cross-Validation and the testing set).
- #4 Split train/val and Cross-Validation(train/val)/test. Pipeline #3 only generate an unbiased testing prediction for a fraction of the dataset (the testing set). To generate an unbiased prediction on every sample of the dataset, pipeline #4 can be used. However, because there is a simple split between the training set and the validation set (one for each Cross-Validation fold), that means that the hyperparameter tuning only leverages a fraction of the dataset. Pipeline #4 is preferred when generating an unbiased testing prediction for every sample is needed but a time-consuming to train to be able to afford a Cross-Validation (RS). Pipeline #4 is therefore the pipeline we used to predict chronological age using medical images or videos, for example.
- #5 Nested Cross-Validation. To fully leverage the whole dataset, a Nested Cross-Validation replaces both the outer split (train/val/test) and the inner splits (train/val) by a Cross-Validation, hence the nesting. It requires the training of a large number of models, and it is therefore usually reserved for models that can be trained relatively quickly. In our case, we used it to train our scalar features-based models (elastic nets, l1/l2 gradient boosted machines, and shallow fully connected neural networks).

Supplementary References

1. Simonyan, K. & Zisserman, A. Very Deep Convolutional Networks for Large-Scale Image Recognition. *arXiv [cs.CV]* (2014).
2. Langner, T., Wikstrom, J., Bjerner, T., Ahlstrom, H. & Kullberg, J. Identifying Morphological Indicators of Aging With Neural Networks on Large-Scale Whole-Body MRI. *IEEE Trans. Med. Imaging* **39**, 1430–1437 (2020).
3. Karargyris, A. *et al.* Age prediction using a large chest x-ray dataset. *Medical Imaging 2019: Computer-Aided Diagnosis* (2019) doi:10.1117/12.2512922.
4. Wang, X. *et al.* Chestx-ray8: Hospital-scale chest x-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases. in *Proceedings of the IEEE conference on computer vision and pattern recognition* 2097–2106 (2017).
5. Wittschieber, D., Vieth, V., Wierer, T., Pfeiffer, H. & Schmeling, A. Cameriere's approach modified for pelvic radiographs: a novel method to assess apophyseal iliac crest ossification for the purpose of forensic age diagnostics. *International Journal of Legal Medicine* vol. 127 825–829 (2013).
6. Martins, D., Vicente, H., Alves, V. & Neves, J. Age Prediction through Pelvis X-Ray Images a Case Based Approach to Problem Solving. *2016 Third International Conference on Mathematics and Computers in Sciences and in Industry (MCSI)* (2016) doi:10.1109/mcsi.2016.058.
7. Wittschieber, D., Vieth, V., Domnick, C., Pfeiffer, H. & Schmeling, A. The iliac crest in forensic age diagnostics: evaluation of the apophyseal ossification in conventional

- radiography. *Int. J. Legal Med.* **127**, 473–479 (2013).
8. Scoles, P. V., Salvagno, R., Villalba, K. & Riew, D. Relationship of iliac crest maturation to skeletal and chronologic age. *J. Pediatr. Orthop.* **8**, 639–644 (1988).
 9. Buckberry, J. L. & Chamberlain, A. T. Age estimation from the auricular surface of the ilium: a revised method. *Am. J. Phys. Anthropol.* **119**, 231–239 (2002).
 10. O'Connor, J. E., Coyle, J., Bogue, C., Spence, L. D. & Last, J. Age prediction formulae from radiographic assessment of skeletal maturation at the knee in an Irish population. *Forensic Sci. Int.* **234**, 188.e1–8 (2014).
 11. Tang, X. *et al.* An Abbreviated Scale for the Assessment of Skeletal Bone Age Using Radiographs of the Knee. *Orthopedics* **41**, e676–e680 (2018).
 12. Fan, F. *et al.* Forensic age estimation of living persons from the knee: Comparison of MRI with radiographs. *Forensic Sci. Int.* **268**, 145–150 (2016).
 13. Dedouit, F. *et al.* Age assessment by magnetic resonance imaging of the knee: a preliminary study. *Forensic Sci. Int.* **217**, 232.e1–7 (2012).
 14. Krämer, J. A. *et al.* Forensic age estimation in living individuals using 3.0T MRI of the distal femur. *International Journal of Legal Medicine* vol. 128 509–514 (2014).
 15. Saint-Martin, P., Rérolle, C., Pucheux, J., Dedouit, F. & Telmon, N. Contribution of distal femur MRI to the determination of the 18-year limit in forensic age estimation. *International journal of legal medicine* vol. 129 619–620 (2015).
 16. O'Connor, J. E., Bogue, C., Spence, L. D. & Last, J. A method to establish the relationship between chronological age and stage of union from radiographic assessment of epiphyseal fusion at the knee: an Irish population study. *Journal of Anatomy* vol. 212 198–209 (2008).
 17. Creo, A. L. & Schwenk, W. F., 2nd. Bone Age: A Handy Tool for Pediatric Providers. *Pediatrics* **140**, (2017).

18. Hao, P. Y. *et al.* Skeletal bone age assessments for young children based on regression convolutional neural networks. *Math. Biosci. Eng.* **16**, 6454–6466 (2019).
19. Westerberg, E. AI-based Age Estimation using X-ray Hand Images : A comparison of Object Detection and Deep Learning models. (2020).
20. Lee, H. *et al.* Fully Automated Deep Learning System for Bone Age Assessment. *J. Digit. Imaging* **30**, 427–441 (2017).
21. Halabi, S. S. *et al.* The RSNA Pediatric Bone Age Machine Learning Challenge. *Radiology* **290**, 498–503 (2019).
22. De Sanctis, V. *et al.* Hand X-ray in pediatric endocrinology: Skeletal age assessment and beyond. *Indian J. Endocrinol. Metab.* **18**, S63–71 (2014).
23. Darmawan, M. F., Yusuf, S. M., Abdul Kadir, M. R. & Haron, H. Age estimation based on bone length using 12 regression models of left hand X-ray images for Asian children below 19 years old. *Leg. Med.* **17**, 71–78 (2015).
24. Schmidt, S. *et al.* Skeletal age determination of the hand: a comparison of methods. *Int. J. Legal Med.* **127**, 691–698 (2013).
25. Chollet, F. Xception: Deep learning with depthwise separable convolutions. in *Proceedings of the IEEE conference on computer vision and pattern recognition* 1251–1258 (2017).
26. Schulze, D. *et al.* Correlation of age and ossification of the medial clavicular epiphysis using computed tomography. *Forensic Sci. Int.* **158**, 184–189 (2006).
27. Zhang, K., Chen, X.-G., Zhao, H., Dong, X.-A. & Deng, Z.-H. Forensic age estimation using thin-slice multidetector CT of the clavicular epiphyses among adolescent western Chinese. *J. Forensic Sci.* **60**, 675–678 (2015).
28. Wittschieber, D. *et al.* Projection radiography of the clavicle: still recommendable for forensic age diagnostics in living individuals? *Int. J. Legal Med.* **129**, 187–193 (2015).

29. Cameriere, R. *et al.* Age estimation in children and young adolescents for forensic purposes using fourth cervical vertebra (C4). *Int. J. Legal Med.* **129**, 347–355 (2015).
30. de Oliveira, F. T. *et al.* Mandibular ramus length as an indicator of chronological age and sex. *Int. J. Legal Med.* **129**, 195–201 (2015).
31. Lewis, A. J. *et al.* Demirjian's method in the estimation of age: A study on human third molars. *J. Forensic Dent. Sci.* **7**, 153–157 (2015).
32. Talabani, R. M., Baban, M. T. & Mahmood, M. A. Age estimation using lower permanent first molars on a panoramic radiograph: A digital image analysis. *J. Forensic Dent. Sci.* **7**, 158–162 (2015).
33. Ge, Z.-P., Ma, R.-H., Li, G., Zhang, J.-Z. & Ma, X.-C. Age estimation based on pulp chamber volume of first molars from cone-beam computed tomography images. *Forensic Sci. Int.* **253**, 133.e1–7 (2015).
34. Van Rossum, G. & Drake, F. L. *The Python Language Reference Manual*. (Network Theory Limited, 2011).
35. Oliphant, T. E. *A guide to NumPy*. vol. 1 (Trelgol Publishing USA, 2006).
36. Walt, S. van der, van der Walt, S., Chris Colbert, S. & Varoquaux, G. The NumPy Array: A Structure for Efficient Numerical Computation. *Computing in Science & Engineering* vol. 13 22–30 (2011).
37. McKinney, W. & Others. Data structures for statistical computing in python. in *Proceedings of the 9th Python in Science Conference* vol. 445 51–56 (Austin, TX, 2010).
38. Hunter, J. D. Matplotlib: A 2D Graphics Environment. *Comput. Sci. Eng.* **9**, 90–95 (2007).
39. Inc, P. T. Collaborative data science. *Montreal: Plotly Technologies Inc Montral* (2015).
40. Clark, A. Pillow Python Imaging Library. *Pillow—Pillow (PIL Fork) 5. 4. 1 documentation* (2018).

41. Virtanen, P. *et al.* SciPy 1.0: fundamental algorithms for scientific computing in Python. *Nature Methods* vol. 17 261–272 (2020).
42. Oliphant, T. E. Python for Scientific Computing. *Computing in Science Engineering* **9**, 10–20 (2007).
43. Millman, K. J., Jarrod Millman, K. & Aivazis, M. Python for Scientists and Engineers. *Computing in Science & Engineering* vol. 13 9–12 (2011).
44. Pedregosa, F. *et al.* Scikit-learn: Machine learning in Python. *the Journal of machine Learning research* **12**, 2825–2830 (2011).
45. Ke, G. *et al.* LightGBM: A Highly Efficient Gradient Boosting Decision Tree. in *Advances in Neural Information Processing Systems 30* (eds. Guyon, I. *et al.*) 3146–3154 (Curran Associates, Inc., 2017).
46. Chen, T. & Guestrin, C. XGBoost: A Scalable Tree Boosting System. in *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* 785–794 (Association for Computing Machinery, 2016).
47. Bergstra, J., Yamins, D. & Cox, D. D. Hyperopt: A python library for optimizing the hyperparameters of machine learning algorithms. in *Proceedings of the 12th Python in science conference* vol. 13 20 (Citeseer, 2013).
48. Abadi, M. *et al.* TensorFlow: Large-scale machine learning on heterogeneous systems. (2015).
49. Chollet, F. & Others. keras. (2015).
50. Kotikalapudi, R. & Others. keras-vis. 2017. URL <https://github.com/raghakot/keras-vis> (2019).
51. Alber, M. *et al.* iNNvestigate neural networks. *J. Mach. Learn. Res.* **20**, 1–8 (2019).
52. Hossain, S., Calloway, C., Lippa, D., Niederhut, D. & Shupe, D. Visualization of

- Bioinformatics Data with Dash Bio. in *Proceedings of the 18th Python in Science Conference* 126–133 (2019).
53. Kohavi, R. & Others. A study of cross-validation and bootstrap for accuracy estimation and model selection. in *Ijcai* vol. 14 1137–1145 (Montreal, Canada, 1995).
 54. Bergstra, J. S., Bardenet, R., Bengio, Y. & Kégl, B. Algorithms for Hyper-Parameter Optimization. in *Advances in Neural Information Processing Systems 24* (eds. Shawe-Taylor, J., Zemel, R. S., Bartlett, P. L., Pereira, F. & Weinberger, K. Q.) 2546–2554 (Curran Associates, Inc., 2011).
 55. Bergstra, J., Yamins, D. & Cox, D. Making a Science of Model Search: Hyperparameter Optimization in Hundreds of Dimensions for Vision Architectures. in (eds. Dasgupta, S. & McAllester, D.) vol. 28 115–123 (PMLR, 2013).
 56. Bergstra, J. & Bengio, Y. Random search for hyper-parameter optimization. *J. Mach. Learn. Res.* **13**, 281–305 (2012).
 57. Kingma, D. P. & Ba, J. Adam: A Method for Stochastic Optimization. *arXiv [cs.LG]* (2014).
 58. Zeiler, M. D. ADADELTA: An Adaptive Learning Rate Method. *arXiv [cs.LG]* (2012).
 59. Hinton, G. Slide 29 of Lecture 6, Geoffrey Hinton coursera’s class.
<http://www.cs.toronto.edu>
http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf.
 60. Nair, V. & Hinton, G. E. Rectified Linear Units Improve Restricted Boltzmann Machines. (2010).
 61. Klambauer, G., Unterthiner, T., Mayr, A. & Hochreiter, S. Self-Normalizing Neural Networks. in *Advances in Neural Information Processing Systems 30* (eds. Guyon, I. et al.) 971–980 (Curran Associates, Inc., 2017).
 62. Prechelt, L. Early Stopping - But When? in *Neural Networks: Tricks of the Trade* (eds. Orr,

- G. B. & Müller, K.-R.) 55–69 (Springer Berlin Heidelberg, 1998).
63. Hochreiter, S. Untersuchungen zu dynamischen neuronalen Netzen. *Diploma, Technische Universität München* **91**, (1991).
 64. Hochreiter, S., Bengio, Y., Frasconi, P., Schmidhuber, J. & Others. Gradient flow in recurrent nets: the difficulty of learning long-term dependencies. (2001).
 65. Alqaraawi, A., Schuessler, M., Weiß, P., Costanza, E. & Berthouze, N. Evaluating saliency map explanations for convolutional neural networks: a user study. in *Proceedings of the 25th International Conference on Intelligent User Interfaces* 275–285 (Association for Computing Machinery, 2020).
 66. Selvaraju, R. R. *et al.* Grad-cam: Visual explanations from deep networks via gradient-based localization. in *Proceedings of the IEEE international conference on computer vision* 618–626 (2017).
 67. Zhou, B., Khosla, A., Lapedriza, A., Oliva, A. & Torralba, A. Learning deep features for discriminative localization. in *Proceedings of the IEEE conference on computer vision and pattern recognition* 2921–2929 (2016).
 68. Wang, Z. & Yang, J. Diabetic Retinopathy Detection via Deep Convolutional Networks for Discriminative Localization and Visual Explanation. *arXiv [cs.CV]* (2017).
 69. Duffy, B. A. *et al.* Regression activation mapping on the cortical surface using graph convolutional networks. (2019).
 70. Le Goallec, A. & Patel, C. J. Age-dependent co-dependency structure of biomarkers in the general population of the United States. *Aging* **11**, 1404–1426 (2019).