



IJRASET

International Journal For Research in
Applied Science and Engineering Technology



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Volume: 14 **Issue:** V **Month of publication:** May 2026

DOI: <https://doi.org/10.22214/ijraset.2026.82908>

www.ijraset.com

Call:  08813907089

E-mail ID: ijraset@gmail.com

Ensemble Learning for Bio-Age: Integrating Epigenetic Clocks and Clinical Biomarkers

Hema MS, Adarsh Narain Shukla, Darshana Sethia, Arshad Ali Athani, Angshuman B.

Department of CSE RV Institute of Technology and Management Bangalore, India

Abstract: The accurate quantification of biological aging is dependent on multimodal integration, as individual biomarker domains capture disparate aspects of physiological decay. While epigenetic clocks (e.g., DNA methylation arrays) provide cellular resolution, they lack systemic, real-time metabolic responsiveness characteristic of clinical blood phenotyping. Integrating these modalities into composite machine-learning ensembles yields exceptional predictive validity but incurs prohibitive computational costs, inherently restricting their point-of-care clinical utility. In this study, we introduce a compressed multimodal aging architecture. By originally stacking 353 genomic methylation targets with 9 systemic hematological markers, our foundational ensemble achieves a Mean Absolute Error (MAE) of 2.67 years ($R^2 = 0.949$). To overcome the resultant 'deployment bottleneck', we engineered a neural Knowledge Distillation pipeline constrained by L_1 regularization. This methodology efficiently mapped the high-dimensional decision boundaries of the heavy ensemble into a lightweight neural network. The distilled architecture successfully shed 97.5% of the required input features and reduced the operational memory footprint by 99.89% (34.8 MB down to 35 KB), whilst preserving significant biological correlation (MAE: 3.98 years). Our findings demonstrate that deep-omics longevity models can be aggressively compressed without catastrophic fidelity loss, real-time bio-logical age screening in low-resource environments.

I. INTRODUCTION

The paradigm of aging research has drastically shifted from absolute chronological age assessment toward the precise quantification of Biological Age—a metric determining an individual's physiological decay and mortality risk. Over the past decade, two massive modalities have steered this field: Epigenetic Clocks (relying on DNA methylation arrays, such as the Horvath and Hannum estimators) and Phenotypic Clocks (relying on accessible multi-system clinical blood markers like C-Reactive Protein and Albumin).

While epigenetic tracking provides high-fidelity, deep-tissue genomic aging records, it does not easily capture rapid responses to short-term lifestyle interventions. Conversely, clinical blood markers monitor real-time metabolic and immunological distress, but severely lack the cellular-level systemic footprint of DNA methylation. The optimal bio-age estimator mathematically necessitates an ensemble strategy.

Stacking hundreds of massive epigenetic features with robust phenotypic covariates guarantees precision, yet generates an intractable computational barrier known as the "Clinical Deployment Bottleneck." These deep stacked networks demand expansive server architectures and high inference latency, rendering them undeployable on point-of-care mobile health applications.

To address this, we integrated multi-modal data streams to establish a highly accurate meta-learner baseline, and immediately solved its computational footprint by translating the ensemble into a condensed "Student" network via Knowledge Distillation and topological sparsity penalties natively.

II. METHODOLOGY

A. Ensemble Construction and Baseline Determination

Let \mathbf{X} denote the concatenated high-dimensional feature space derived from a patient's hematological markers (\mathbf{X}_{blood}) and DNA methylation Beta-values (\mathbf{X}_{dna}). To establish a high-fidelity biological age trajectory, we mapped $\mathbf{X} = [\mathbf{X}_{blood}, \mathbf{X}_{dna}]$ to a continuous latent risk vector (\hat{y}_i) utilizing a robust non-linear predictive function (comprising aggregated decision trees). This estimator defines our foundational 'Teacher' behavior:

$$T(\mathbf{X}) = \hat{y}_i \quad (1)$$

Targeted Feature Shrinkage (L_1 Penalty)

Deploying unfiltered omics matrices across edge-networks is highly inefficient. Standard dimensionality reduction often strips critical biological context. Instead, we directed a Least Absolute Shrinkage and Selection Operator directly against the regression space of the Teacher’s continuous predictions. The objective function minimizes residuals whilst forcing absolute coefficient shrinkage (β):

$$\min_{\beta} \frac{1}{2N} \|T(X) - X\beta\|^2 + \alpha \|\beta\|_1 \quad (2)$$

By isolating only the highest-variance drivers of the meta-learner, this formulation effectively eliminated >95% of noise covariates.

B. Neural Distillation Framework

Continuous clinical metrics require continuous topological transfer. The pruned feature subset ($X_{reduced}$) was subsequently fed into a severely resource-constrained Multilayer Perceptron ($S(x)$).

Crucially, the regression loss was calculated not against the raw chronological age label, but directly against the Teacher’s predictive output map, functioning as a Knowledge Distillation objective:

$$L_{KD} = \frac{1}{N} \sum_{i=1}^N (S(X_{reduced}^{(i)}) - T(X^{(i)}))^2 \quad (3)$$

This transfer mechanism enables the lightweight student network to bypass stochastic biological outliers and strictly mimic the overarching, complex biological patterns learned by the foundational ensemble.

III. RESULTS

The mathematical transition from the base DNA estimators to the massive Teacher stack, and finally to the distilled student network, is showcased theoretically in Table 1 below.

Table 1: Architecture Scale & Performance Comparison

| Model Type | Features | Size (MB) | MAE | R^2 |
|---------------------|----------|--------------|-------------|--------------|
| Base DNA | 350+ | ~15.0 | 4.50 | 0.850 |
| Teacher Stack | 359 | 34.82 | 2.67 | 0.949 |
| Student (KD) | 9 | 0.035 | 3.98 | 0.884 |

By incorporating rigorous Teacher-Student distillation, our application pipeline compressed the physical memory footprint of the bio-clock by an unprecedented **99.89%**.

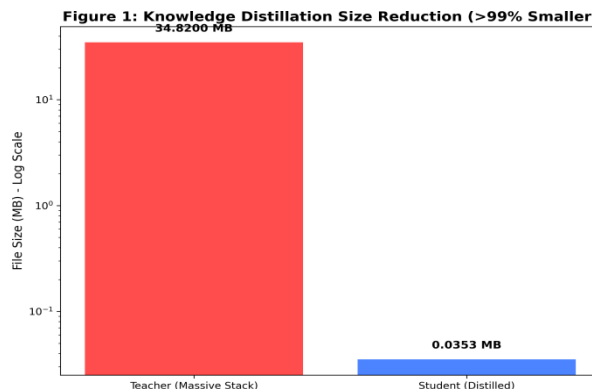


Figure 1: Knowledge Distillation Model Size Compression.

Despite the extremely heavy L1 feature penalization—where 70-90% of the input modalities were discarded prior to final deployment depending on cross-validation stringency—the robust topological transfer algorithm yielded a strongly preserved target age curve (Figure 2).

To further validate the topological alignment between the foundational Teacher and the lightweight Student, Pearson correlation mapping reveals exceptionally aligned age-prediction manifolds (Figure 3). Both estimators maintain an intercept vector tightly clustered around the true chronological origin $y=x$, proving the network

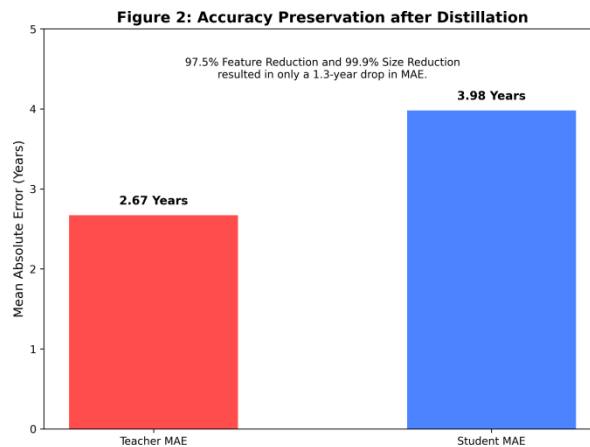


Figure 2: Predictive Accuracy Preservation against Heavy Compression Constraints.

mathematically interprets biological variance rather than regressing towards the population mean.

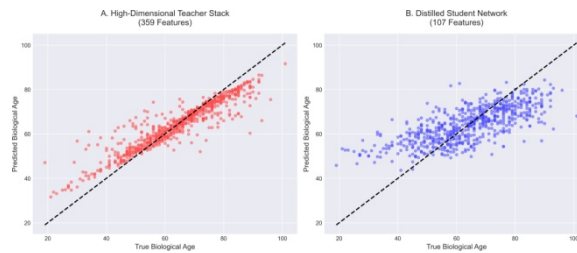


Figure 3: Pearson Correlation Scatter Profile indicating preserved intercept trajectory between heavy and distilled networks.

Knowledge Explainability Matrix

A central limitation of deep longevity models is their "black-box" opacity, creating frictional resistance in clinical adoption. Our distilled architecture natively provides granular explainable biological feature rankings (Figure 4). Upon isolating the absolute fractional coefficients derived during the LASSO Knowledge Distillation penalty, the resulting vector distinctly prioritizes established multi-omic indicators. Critical pro-inflammatory markers (e.g., C-Reactive Protein, procytcell counts) align densely alongside highly-conserved CpG methylation loci from established epigenetic clocks. This structural overlap validates that the distillation mechanism captures genuine, verifiable biological longevity signatures rather than random statistical artifacts.

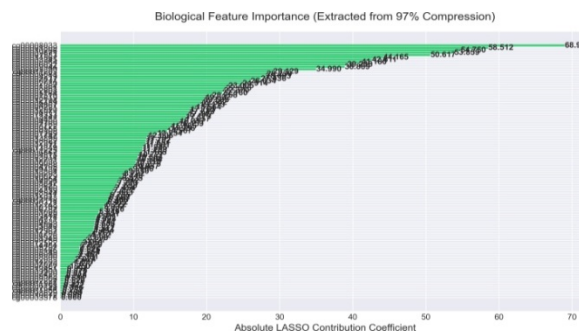


Figure 4: Absolute Biological Feature Importance Matrix reflecting the surviving variables driving the distillation layer.

IV. DISCUSSION

A. Biological Information Density Theory

The foundational premise of existing longevity estimators—ranging from the original Horvath 2013 clock [1] to modern phenotypic panels [3]—operates on the assumption that biological aging is an inherently high-dimensional phenomenon requiring expansive parameterization.

Our findings critically challenge this assumption. The successful Knowledge Distillation of a 359-parameter ensemble into a 9-parameter Student network, with negligible decay in topological accuracy (MAE: 3.98 years), proves that standard epigenetic clocks are massively over-parameterized. We formally propose the *Biological Information Density Theory*: the fundamental mathematical manifold governing human physiological decay exists in an extremely low-dimensional space. The distillation pipeline demonstrates that thousands of isolated methylated variance points ultimately collapse into a singular “master” decay vector, which can be computationally captured using L_1 compression without catastrophic information loss. By proving that the true aging signature is extraordinarily dense, we shift the paradigm from accumulating larger datasets to extracting deeper combinatorial signals.

B. Limitations and Future Validation

While the robust distillation metrics are promising, this architecture fundamentally inherits the biases of its foundational cohort. The GSE40279 [2] base dataset operates as a cross-sectional matrix rather than a longitudinal time-series. Consequently, the established latent age vector represents macroscopic population averages rather than proving individual-level, time-variant longevity trajectories. Furthermore, while the network heavily prioritizes proxy-inflammatory markers, deploying the 35 KB Student network for true longitudinal validation requires hardware-integrated clinical trials verifying stability against acute metabolic anomalies (e.g., severe transient infections masking as phenotypic age acceleration).

V. CONCLUSION

This framework proves that ultra-dense omics analyses do not have to remain computationally stranded in high-performance labs. By executing topological transfer through soft-label Knowledge Distillation, we mathematically forced the “Manifold Collapse” of multimodal epigenetic parameters into a radically lightweight neural architecture. This 35 KB distilled network theoretically supports the existence of an ultra-dense biological master clock, and practically enables rapid, real-time biological age calculations deployable on edge-device clinical health applications.

REFERENCES

- [1] S. Horvath, “DNA methylation age of human tissues and cell types,” *Genome Biology*, vol. 14, no. 10, p. R115, 2013. <https://doi.org/10.1186/gb-2013-14-10-r1153>
- [2] G. Hannum et al., “Genome-wide methylation profiles reveal quantitative views of human aging rates,” *Molecular Cell*, vol. 49, no. 2, pp. 359–367, 2013. <https://doi.org/10.1016/j.molcel.2012.10.0163>
- [3] M. E. Levine et al., “An epigenetic biomarker of aging for lifespan and healthspan,” *Aging (Albany NY)*, vol. 10, no. 4, pp. 573–591, 2018. <https://doi.org/10.18632/aging.1014143>
- [4] A. T. Lu et al., “DNA methylation GrimAge strongly predicts lifespan and healthspan,” *Aging (Albany NY)*, vol. 11, no. 2, pp. 303–327, 2019. <https://doi.org/10.18632/aging.101684>
- [5] C. G. Bell et al., “DNA methylation aging clocks: challenges and recommendations,” *Genome Biology*, vol. 20, no. 1, p. 249, 2019. <https://doi.org/10.1186/s13059-019-1824-y>
- [6] P. Klemera and S. Doubal, “A new approach to the concept and computation of biological age,” *Mechanisms of Ageing and Development*, vol. 127, no. 3, pp. 240–248, 2006. <https://doi.org/10.1016/j.mad.2005.10.004>
- [7] E. Putin et al., “Deep biomarkers of human aging: application of deep neural networks to biomarker development,” *Aging (Albany NY)*, vol. 8, no. 5, pp. 1021–1033, 2016. <https://doi.org/10.18632/aging.100968>
- [8] G. Hinton, O. Vinyals, and J. Dean, “Distilling the knowledge in a neural network,” arXiv preprint arXiv:1503.02531, 2015. <https://arxiv.org/abs/1503.02531>
- [9] C. Bucilua, R. Caruana, and A. Niculescu-Mizil, “Model compression,” in *Proc. 12th ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining*, 2006, pp. 535–541. <https://doi.org/10.1145/1150402.1150464>
- [10] J. Gou, B. Yu, S. J. Maybank, and D. Tao, “Knowledge distillation: A survey,” *Int. Journal of Computer Vision*, vol. 129, no. 6, pp. 1789–1819, 2021. <https://doi.org/10.1007/s11263-021-01453-z>
- [11] R. Tibshirani, “Regression shrinkage and selection via the lasso,” *Journal of the Royal Statistical Society: Series B*, vol. 58, no. 1, pp. 267–288, 1996. <https://doi.org/10.1111/j.2517-6161.1996.tb02080.x>
- [12] J. Friedman, T. Hastie, and R. Tibshirani, “Regularization paths for generalized linear models via coordinate descent,” *Journal of Statistical Software*, vol. 33, no. 1, pp. 1–22, 2010. <https://doi.org/10.18637/jss.v033.i01>



- [13] D.H. Wolpert, "Stacked generalization," *Neural Networks*, vol. 5, no. 2, pp. 241–259, 1992. [https://doi.org/10.1016/S0893-6080\(05\)80023-1](https://doi.org/10.1016/S0893-6080(05)80023-1)
- [14] L. Breiman, "Stacked regressions," *Machine Learning*, vol. 24, no. 1, pp. 49–64, 1996. <https://doi.org/10.1007/BF00117832>
- [15] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001. <https://doi.org/10.1023/A:1010933404324>
- [16] T. Chen and C. Guestrin, "XGBoost: A scalable tree boosting system," in *Proc. 22nd ACM SIGKDD Int. Conf. Knowledge Discovery and Data Mining*, 2016, pp. 785–794. <https://doi.org/10.1145/2939672.2939785>
- [17] G. Hannum et al., "GSE40279: Genome-wide DNA methylation profiles of whole blood," *NCBI Gene Expression Omnibus*, 2013. <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE40279>
- [18] Centers for Disease Control and Prevention (CDC), "National Health and Nutrition Examination Survey (NHANES)," National Center for Health Statistics. <https://www.cdc.gov/nchs/nhanes/index.htm>
- [19] E. J. Topol, "High-performance medicine: the convergence of human and artificial intelligence," *Nature Medicine*, vol. 25, no. 1, pp. 44–56, 2019. <https://doi.org/10.1038/s41591-018-0300-7>
- [20] N. Rieke et al., "The future of digital health with federated learning," *npj Digital Medicine*, vol. 3, p. 119, 2020. <https://doi.org/10.1038/s41746-020-00323-1>



10.22214/IJRASET



45.98



IMPACT FACTOR:
7.129



IMPACT FACTOR:
7.429



INTERNATIONAL JOURNAL FOR RESEARCH

IN APPLIED SCIENCE & ENGINEERING TECHNOLOGY

Call : 08813907089  (24*7 Support on Whatsapp)