

Using Deep Learning Model to Predict Individual Health Outcomes Due to Migration

Naaisha Agarwal
Advisor: Dr. Kevin Crowthers

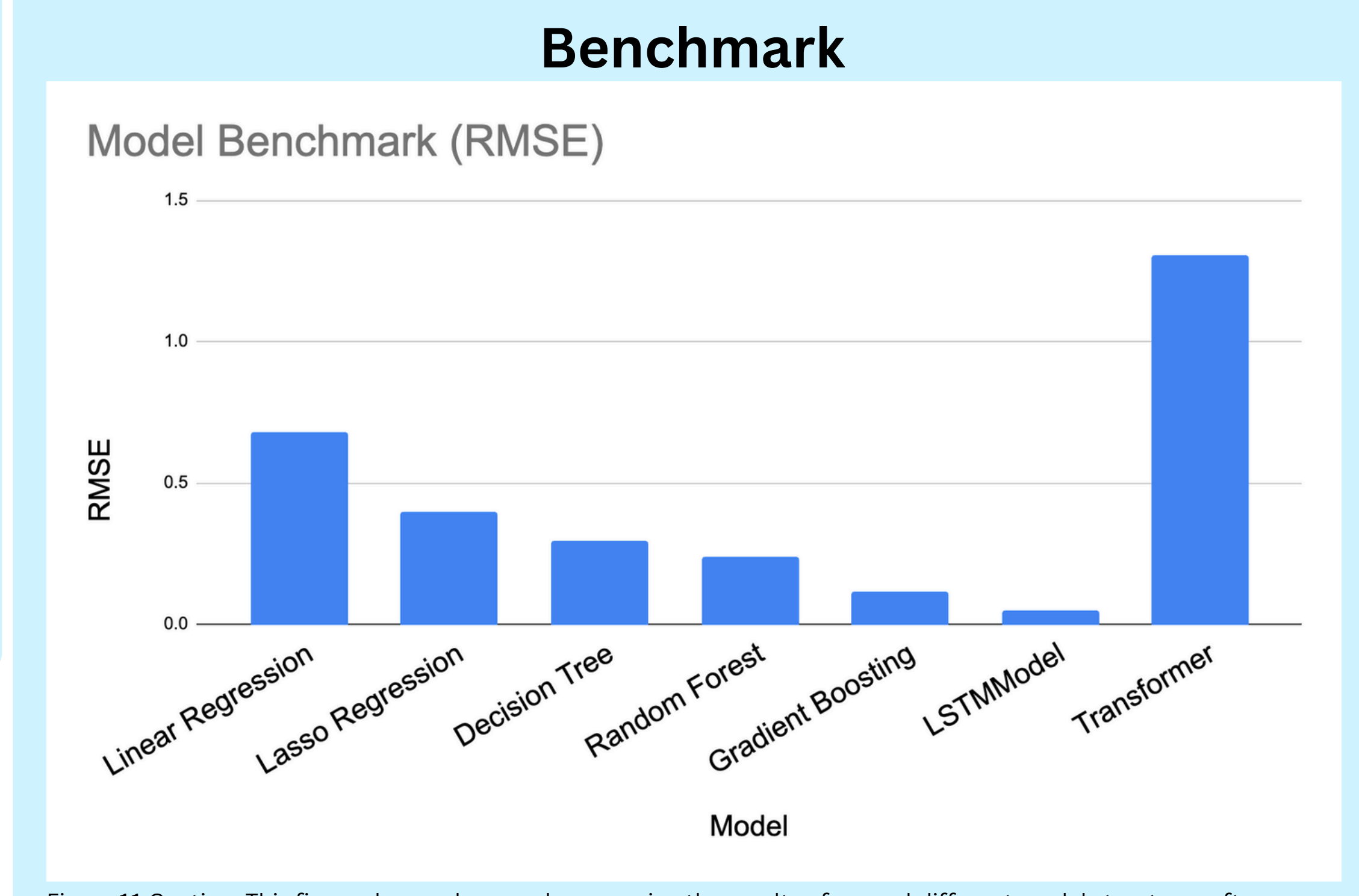
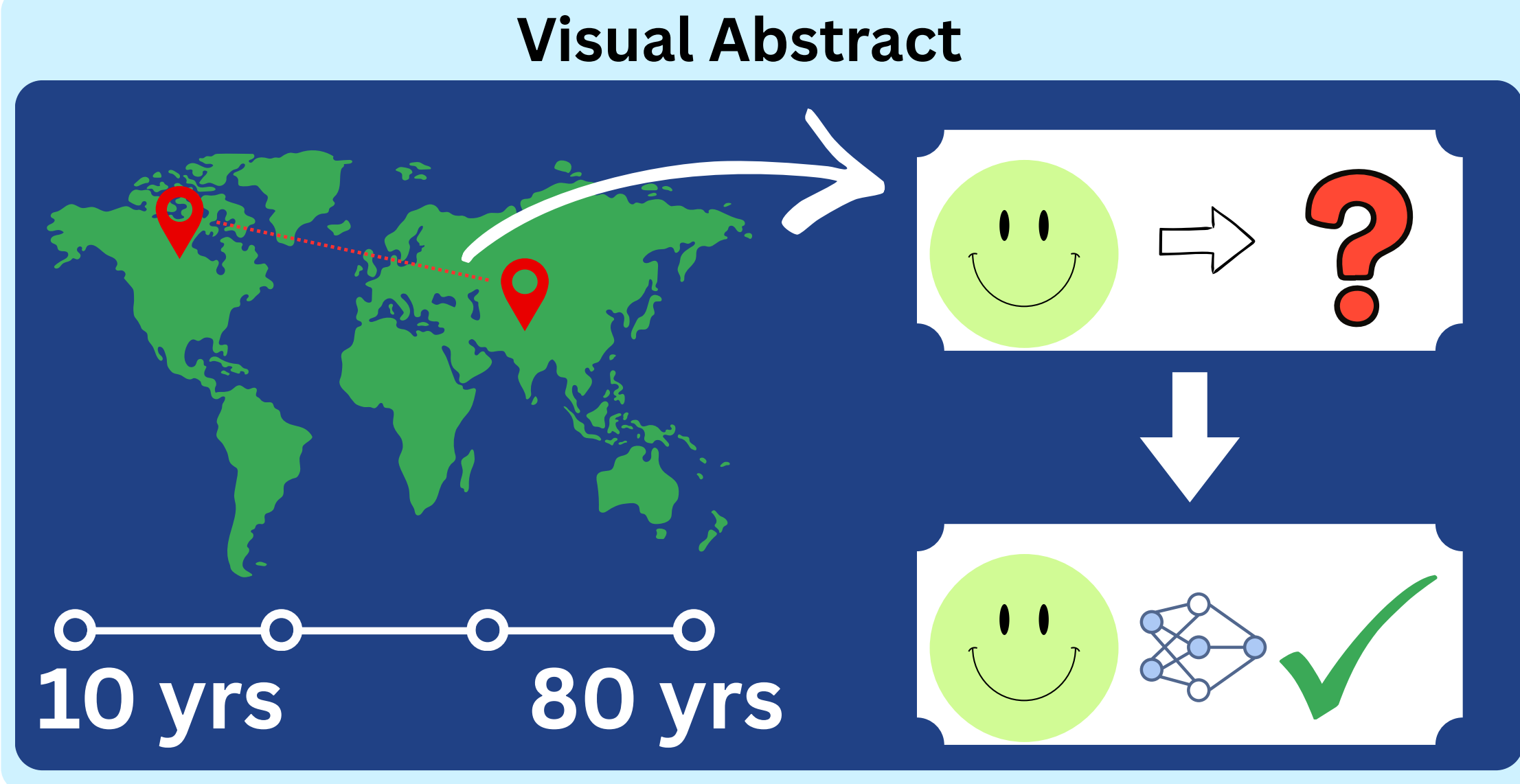
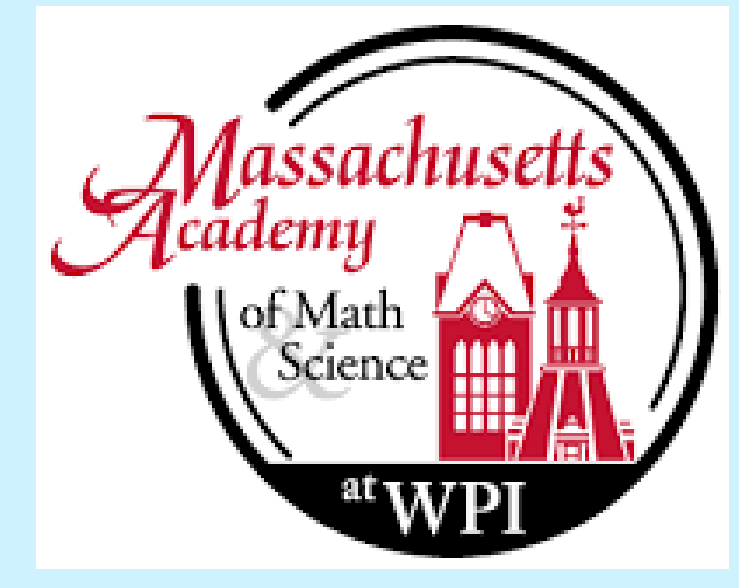


Figure 11 Caption: This figure shows a bar graph comparing the results of several different model structures after training on the curated dataset, representing a benchmark for this problem. RMSE was used as the error metric for comparison. LSTMModel, the model trained for this project, has the lowest RMSE, demonstrating that it has the least loss compared to all other model structures, thereby performing the strongest.

Background

Problem

40 million migrants each year

- Migrants initially experience better resilience to diseases, but later on experience worsened health (Paez-Deggeller, 2025)
- All migrant health studies so far are **population level** studies – does not help individual understand their specific situation

Opportunity

- Recent shift toward **personalized treatments** and risk diagnosis → provides increased accuracy for patients and doctors (Johnson et al., 2021)
- Deep learning models** have been used to predict disease risk using **Electronic Health Record (EHR)** data through **temporal patterns** (Amirahmadi et al., 2023)

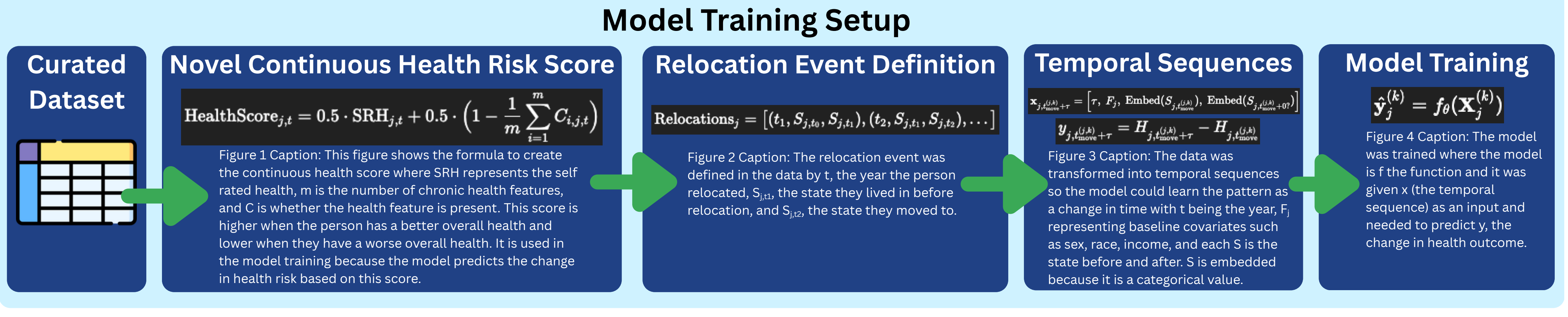
“Healthy migrant Effect”

Engineering Problem:

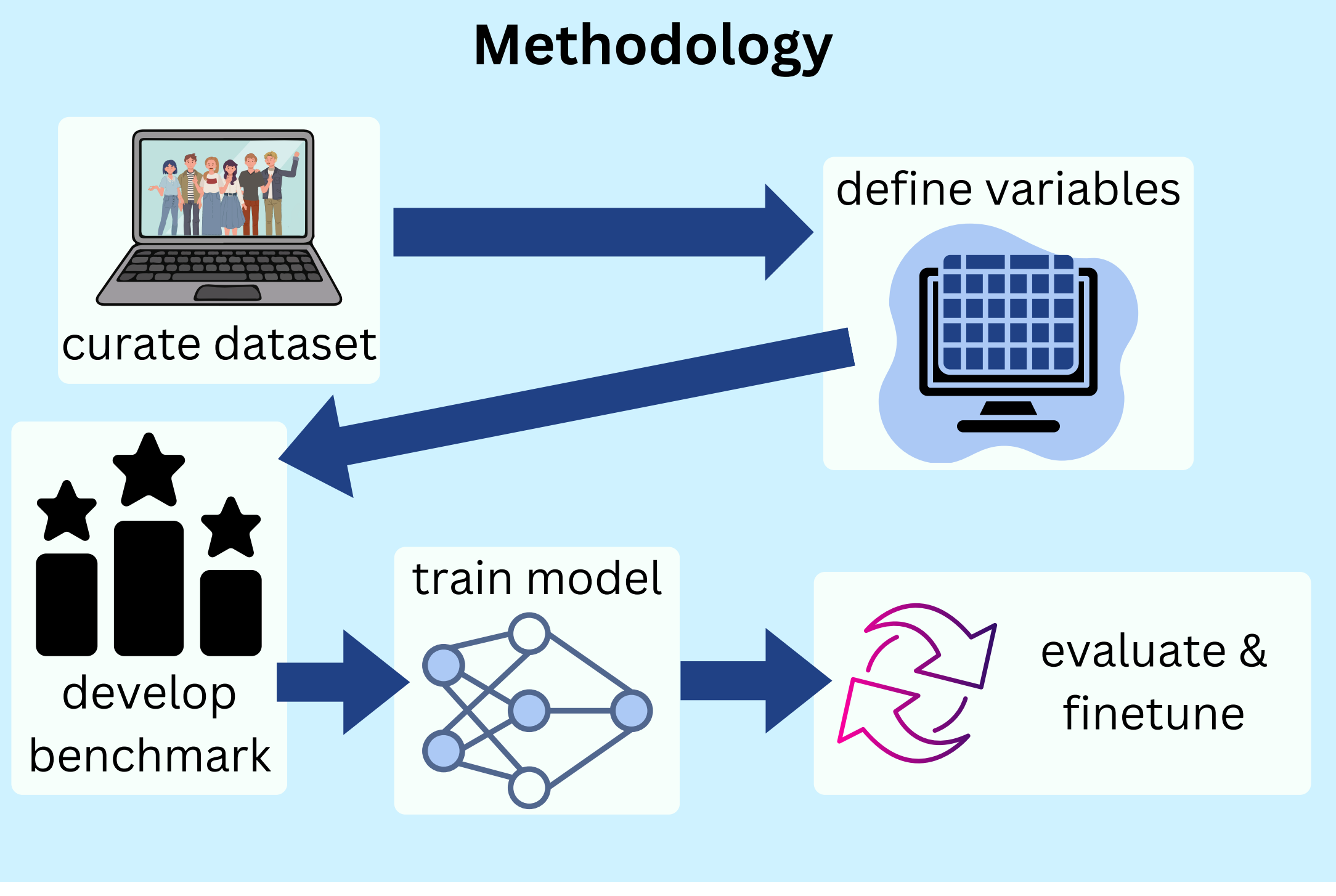
Migrants face several **health changes** after **relocating**, yet they don't have a way to **predict** their **individual** health changes so they can be **prepared**.

Engineering Goal

The goal is to engineer a **model** that **accurately** identifies **potential health risks** individuals will face based on their personal demographics, history, and location they are moving to.



- ### Key Challenges
- Dataset** needs to be: large, long time period, migrant only, baseline and resulting health, individual demographics, residential history
 - Modeling the **causal** and **temporal** nature of the data for the highest accuracy is tricky
 - Model needs to be **interpretable** for the user



- ### Analysis
- LSTM performs with the least error compared to several other model structures justifying its use as the model structure for this project
 - Model has a consistent decreasing loss across epochs showing its need for 200 epochs of training
 - Model's error distribution is balanced, showing that it does not bias toward either direction
 - Model's predictions generally match trajectory of health changes, although it struggles with predicting extreme risk accurately and still shows variance

- ### Conclusion
- First** individualized migrant health prediction tool
 - Empowers **informed decision-making** for migrants
 - Establishes a **scalable framework** for personalized migrant health predictions and preparedness
 - Can be further used by doctors and healthcare professionals to improve suggestions for individuals based on relocation history
 - Shows model's ability to learn temporal and causal data

- ### Key Contributions
- A new **dataset** containing demographics, baseline health, health after relocation, and residential histories specifically for migrants
 - A **benchmark** various model structures on their ability to train to this dataset type
 - A novel continuous health risk score helping define health risk over years
 - A new **model** that is trained on the dataset and accurately predicts health outcomes for people who live in one area and migrate to another

Decision Matrix

Criterion	Weight	Linear Models	Deep Learning	Tree-Based Models
Temporal Modeling	25%	0.1	0.983	0.2
Feature Representation	20%	0.3	1	0.65
Overfitting Resistance	15%	0.7	0.645	0.85
Training Speed	10%	1	0.55	0.95
Interpretability	10%	0.9	0.325	0.9
Handling Missing / Sparse Data	10%	0.3	0.525	0.95
Scalability to Large Datasets	10%	0.6	0.917	0.8
Weighted Score	100%	0.47	0.77	0.67

About:

Sourced from Panel Study of Income Dynamics (Goukova et al., 2026)

Key Facts:

- Total participants: 18890
- 879 features tracked
- 6 health features

~17 million data points

Figure 9 Caption: This shows the distribution of the families based on their starting state in 1999. Although certain states have many more participants than other states, overall the dataset shows a balanced mix of participants from many states, showing how it is more representative of the entire United States and not just one state or group of states.

Figure 10 Caption: This figure shows a correlation heatmap of the highest correlated features in the dataset. Almost all the features are wealth and income across many years. This makes sense because one feature across years should show correlation. This shows that the dataset serves its function as it should show a correlation between the same variable across years.

- ### Future Work
- Map each state to environmental factors
 - Create counterfactual experiments for comparisons
 - Test model on other datasets to evaluate performance
 - Apply model to specific disease predictions to improve predictions based on relocations and suggest potential relocations for improved disease trajectories

References

Amirahmadi, A., Ohlsson, M., & Etmnani, K. (2023). Deep learning prediction models based on EHR trajectories: A systematic review. *Journal of Biomedical Informatics*, 144, 104430. <https://doi.org/10.1016/j.jbi.2023.104430>

Chi, G., Abel, G. J., Johnston, D., Giraudy, E., & Bailey, M. (n.d.). Measuring global migration flows using online data. *Proceedings of the National Academy of Sciences of the United States of America*, 122(18), e2409418122. <https://doi.org/10.1073/pnas.2409418122>

Goukova, E., Andreski, P., & Schoeni, R. F. (2026). Panel Study of Income Dynamics.

Johnson, K. B., Wei, W., Wiseratane, D., Frisse, M. E., Misulis, K., Rhee, K., Zhao, J., & Snowden, J. L. (2021). Precision Medicine, AI, and the Future of Personalized Health Care. *Clinical and Translational Sciences*, 14(1), 86–93. <https://doi.org/10.1111/cts.12884>

Paez-Deggeller, V. (2025, August 25). Top Statistics on Global Migration and Migrants. <https://www.migrationpolicy.org/article/top-statistics-global-migration-migrants>

Team, G., Georgiev, P., Lei, V. I., Burnell, R., Bai, L., Gulati, A., Tanzer, G., Vincent, D., Pan, Z., Wang, S., Mariooryad, S., Ding, Y., Geng, X., Alcober, F., Frostig, R., Omernick, M., Walker, L., Paduraru, C., Sorokin, C., ... Vinyals, O. (2024). Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context (No. arXiv:2403.05530). arXiv. <https://doi.org/10.48550/arXiv.2403.05530>