

A Machine-Learned Comorbidity Index

Suleman Baloch¹, Kishlay Jha², Alberto M. Segre¹, Philip M. Polgreen³, Bijaya Adhikari¹

University of Iowa

¹Department of Computer Science, University of Iowa

²Department of Electrical and Computer Engineering, University of Iowa

³Department of Internal Medicine, University of Iowa



Connect
with me!



Background and Motivation

A **comorbidity index** summarizes the diagnosis burden recorded for a hospital admission into a single scalar score.

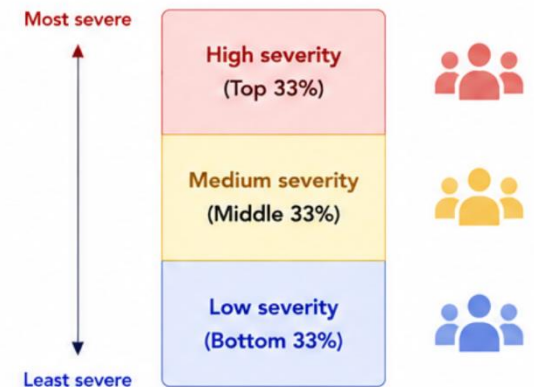
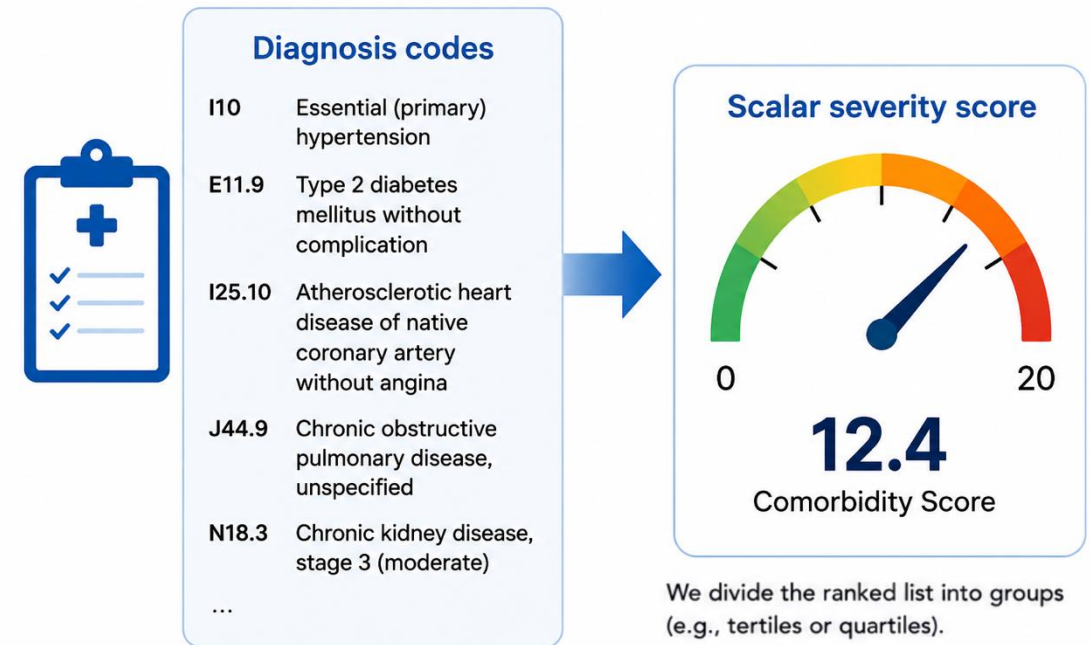
For a hospital admission:

Diagnosis codes → **scalar severity score**

Higher score usually means the admission has more severe or complex comorbid illness.

Comorbidity scores are widely used in healthcare:

1. Patient stratification
2. Risk adjustment
3. Comparing outcomes across patient groups
4. Identifying high-risk admissions






















Traditional Comorbidity Scores

Widely used traditional comorbidity indices include:

1. Charlson Comorbidity Index (CCI¹)
2. van Walraven Elixhauser Index (ECI^{2,3})

How they work:

1. Identify predefined comorbidity conditions from ICD codes
2. Assign each condition a fixed weight
3. Sum the weights to get one scalar score

	Comorbid Condition	Weight
1	 Myocardial infarction	1
2	 Congestive heart failure	1
3	 Peripheral vascular disease	1
4	 Cerebrovascular disease	1
5	 Dementia	1
6	 Chronic pulmonary disease	1
7	 Connective tissue disease / rheumatologic disease	1
8	 Peptic ulcer disease	1
9	 Mild liver disease	1
10	 Diabetes without end-organ damage	1
11	 Diabetes with end-organ damage	2
12	 Hemiplegia or paraplegia	2
13	 Moderate or severe renal disease	2
14	 Any tumor / malignancy without metastasis	2
15	 Leukemia	2
16	 Lymphoma	2
17	 Moderate or severe liver disease	3
18	 Metastatic solid tumor	6
19	 AIDS / HIV	6

Charlson Comorbidity Index (CCI) Formula

$$CCI = \sum_k w_k \times I\{\text{condition } k \text{ present}\}$$

Where:

- w_k = weight for comorbid condition k (from table above)
- $I\{\text{condition } k \text{ present}\} = 1$ if condition k is present, 0 otherwise

1. Charlson ME, Pompei P, Ales KL, MacKenzie CR. A new method of classifying prognostic comorbidity in longitudinal studies: development and validation. J Chronic Dis. 1987;40(5):373-83. doi: 10.1016/0021-9681(87)90171-8. PMID: 3558716.
2. Elixhauser A, Steiner C, Harris DR, Coffey RM. Comorbidity measures for use with administrative data. Med Care. 1998 Jan;36(1):8-27. doi: 10.1097/00005650-199801000-00004. PMID: 9431328.
3. van Walraven C, Austin PC, Jennings A, Quan H, Forster AJ. A modification of the Elixhauser comorbidity measures into a point system for hospital death using administrative data. Med Care. 2009 Jun;47(6):626-33. doi: 10.1097/MLR.0b013e31819432e5. PMID: 19433995.

Limitations of Comorbidity Indices

Limitation 1: Mortality-centric fixed weights

Traditional indices were designed mainly for mortality but are often used for other outcomes.

Comorbidity	Score
Prior myocardial infarction	1
Congestive heart failure	1
Peripheral vascular disease	1
Cerebrovascular disease	1
Dementia	1
Chronic pulmonary disease	1
Rheumatologic disease	1
Peptic ulcer disease	1
Mild liver disease	1
Diabetes	1
Cerebrovascular (hemiplegia) event	2
Moderate-to-severe renal disease	2
Diabetes with chronic complications	2
Cancer without metastases	2
Leukemia	2
Lymphoma	2
Moderate or severe liver disease	3
Metastatic solid tumor	6
Acquired immuno-deficiency syndrome (AIDS)	6

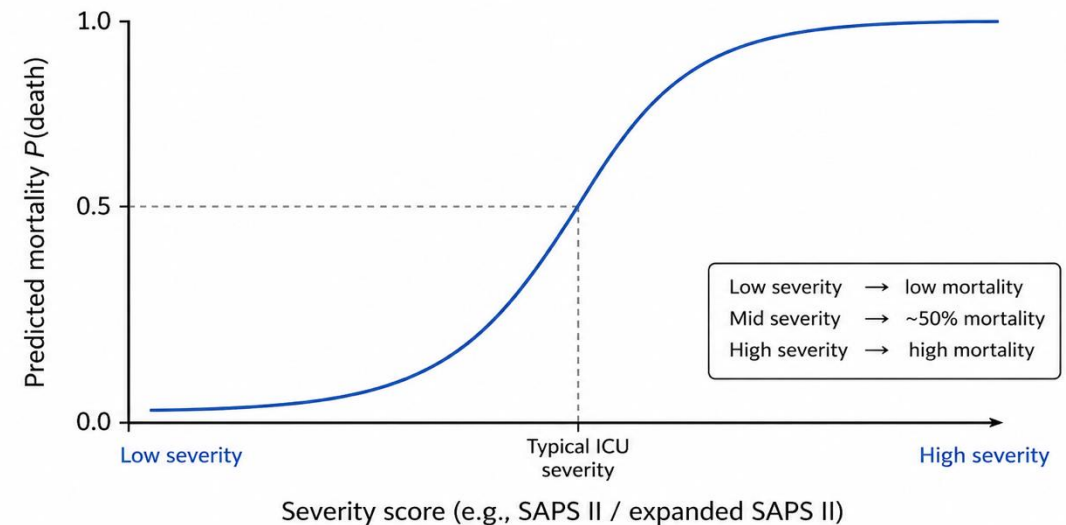
Fixed weights calibrated to predict mortality

Limitation 2: Linear, rule-based structure

They are linear and rule-based, so they may miss nonlinear diagnosis–outcome relationships.

Mortality model: logistic link

$$P(\text{death}) = \frac{e^{\text{logit}}}{1 + e^{\text{logit}}}$$



Framing Comorbidity Scoring as a Machine Learning Problem

These two limitations suggest a **learning problem**:

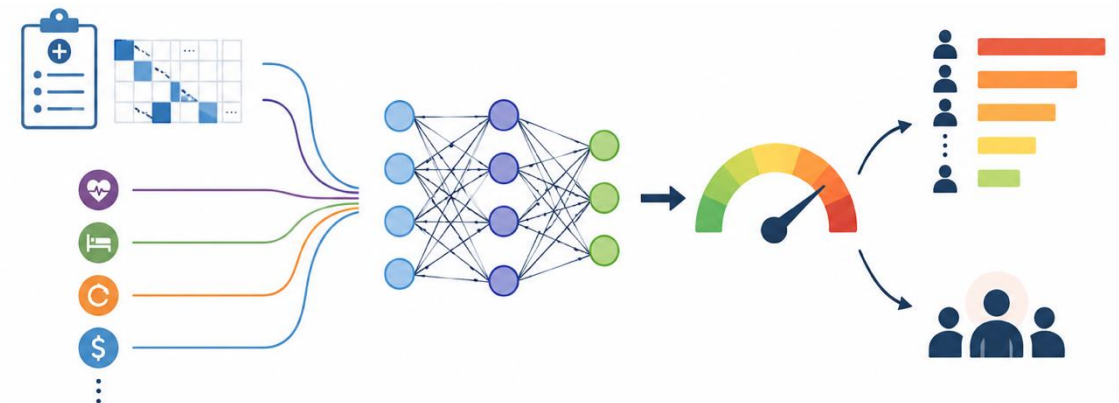
“Can we learn a scalar comorbidity score from data that aligns with multiple outcomes and captures nonlinear risk–outcome interactions?”

We therefore formulate comorbidity scoring as a **machine learning problem**: learn a score that is

1. **Scalar**: preserves the practical simplicity of traditional indices
2. **Multi-outcome aligned**: captures shared comorbidity burden across outcomes
3. **Nonlinear**: captures complex diagnosis–outcome relationships
4. **Stratifiable**: ranks admissions and identifies a high-severity group
5. **Data-driven**: learned from diagnosis codes and outcomes

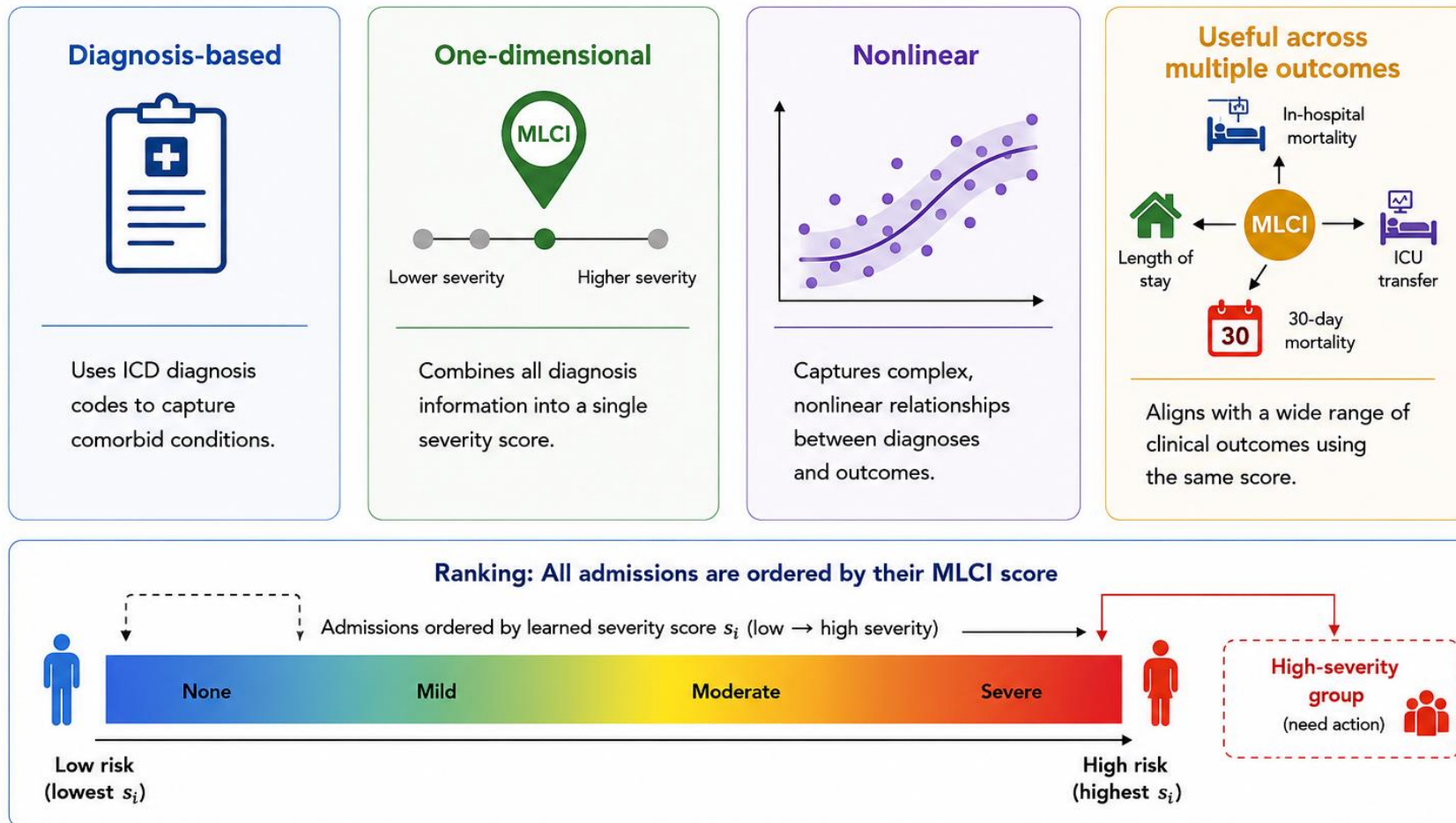
To address this, we introduce **MLCI**:

A MACHINE-LEARNED COMORBIDITY INDEX



MLCI Overview

MLCI learns one data-driven scalar admission-level score by maximizing normalized dependence with multiple clinical outcomes.



MLCI: Intuition

Severity-ordering intuition: Admissions with greater latent comorbidity burden should receive higher learned scores after orientation. Therefore, after orientation, higher MLCI score = higher learned severity.

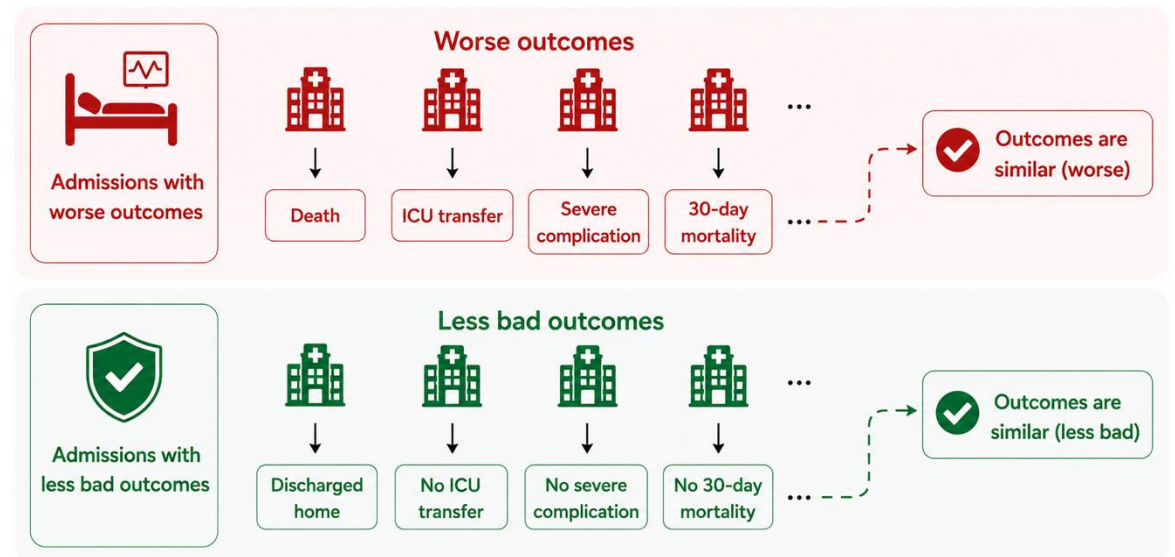
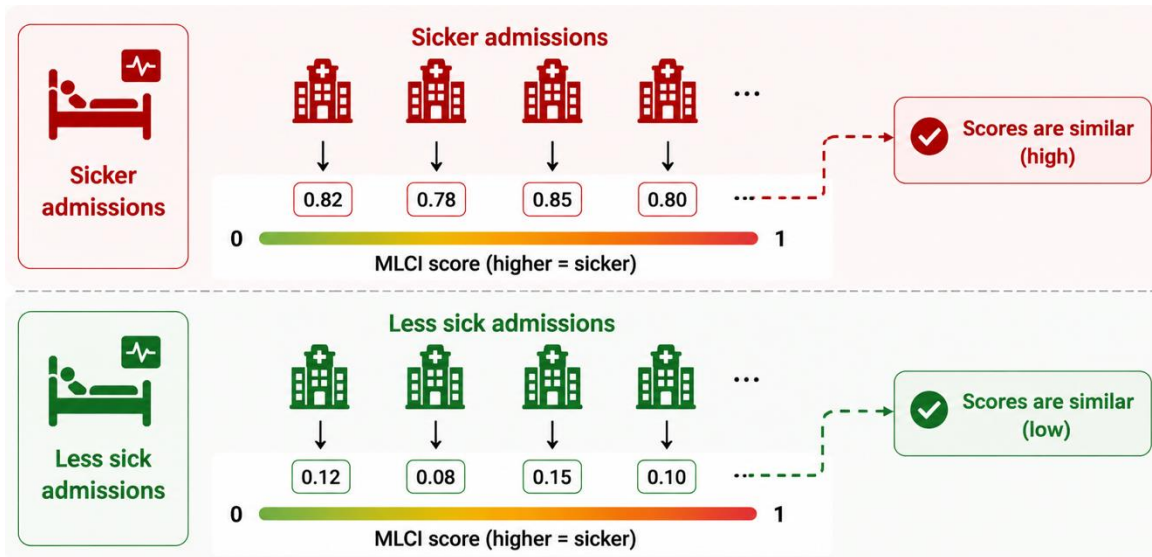
Intuition: For every pair of admissions, we pose two questions

Are their scores similar?

→ Sicker admissions should have similar high scores; less sick admissions should have similar low scores.

Are their outcomes similar?

→ Admissions with matching outcome labels should be similar.



MLCI: Intuition

Pairwise view: For every pair of admissions, MLCI compares two similarities:

- Score similarity:** Are their MLCI scores close?
- Outcome similarity:** Are their clinical outcome labels similar?

MLCI then maximizes the alignment between these two similarity matrices using **nHSIC**.

Goal: admissions that are close in MLCI score should also have similar outcome labels.

For every pair of admissions (i, j) , compare:

Score similarity : Are their MLCI scores similar?

Outcome similarity : Are their outcome labels similar?

Admissions	MLCI score (higher = sicker)
A ₁	0.82
A ₂	0.65
A ₃	0.30
A ₄	0.10

1) Score Similarity Matrix
(from MLCI scores)

	A ₁	A ₂	A ₃	A ₄
A ₁	1.00	0.80	0.20	0.05
A ₂	0.80	1.00	0.30	0.10
A ₃	0.20	0.30	1.00	0.70
A ₄	0.05	0.10	0.70	1.00

Low similarity High similarity
(1 = very similar, 0 = very dissimilar)

2) Outcome Similarity Matrix
(from outcome labels)

	A ₁	A ₂	A ₃	A ₄
Mortality	1	1	0	0
ICU Transfer	1	1	0	0
30-day Mortality	1	0	1	0
	A ₁	A ₂	A ₃	A ₄
A ₁	1	0	1	0
A ₂	0	1	0	0
A ₃	1	0	1	1
A ₄	0	0	1	1

1 Similar (labels match) **0** Not similar (labels differ)

(Computed the same way for all outcomes, then combined across outcomes in nHSIC)

MLCI Architecture

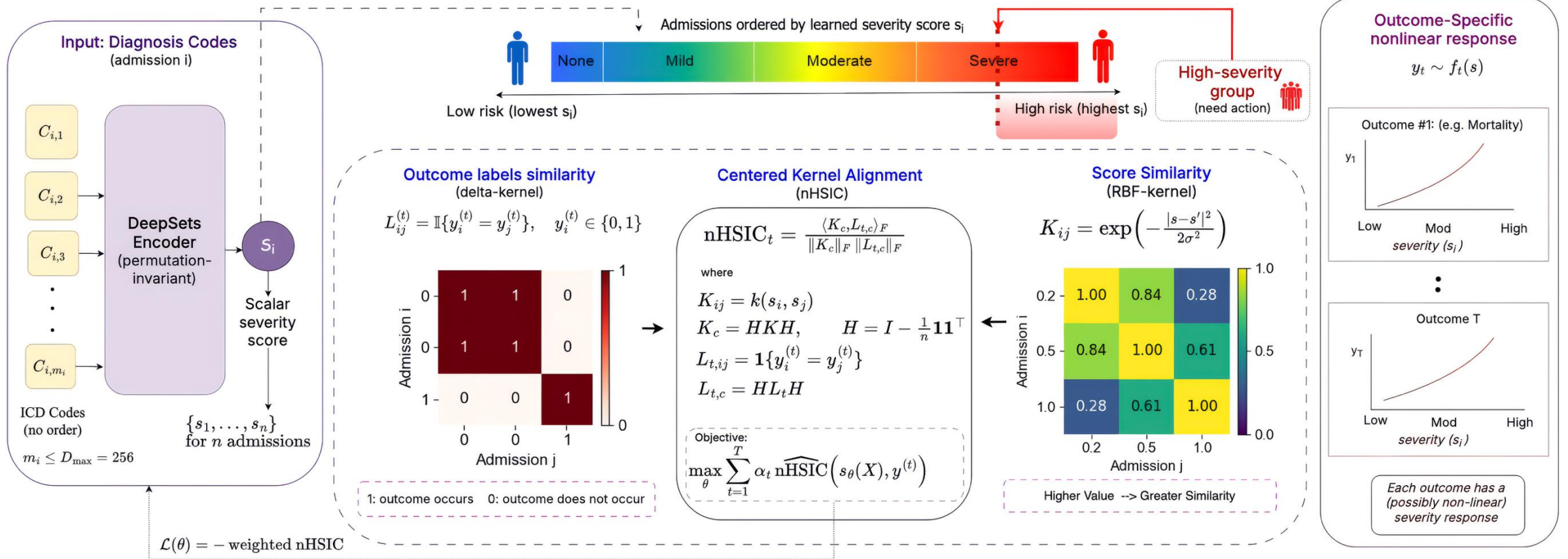


Figure guide: A DeepSets encoder maps diagnosis codes to a scalar MLCI score. During training, MLCI learns this score by aligning score similarity with outcome-label similarity using nHSIC. At inference, only diagnosis codes are needed to compute the MLCI score and rank admissions. Post hoc outcome-specific curves map the learned MLCI score to outcome-specific risk patterns.

Theoretical Analysis: Oracle Setup

Theoretical goal: Our analysis characterizes when one oracle comorbidity score can support an approximate shared admission-level ordering across multiple clinical outcomes.

Setup: We analyze this in a finite-sample oracle setting.

1. Assume each admission has an unobserved latent severity

$$z_i = \text{latent disease burden for admission } i$$

2. Each outcome depends on severity through its own function

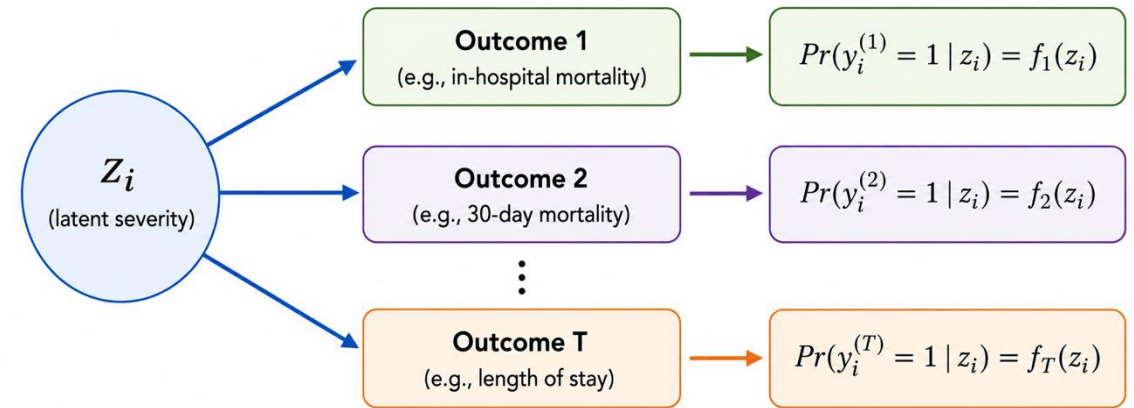
$$\Pr(y_i^{(t)} = 1 \mid z_i) = f_t(z_i)$$

Thus, different outcomes may respond differently to severity.

Goal: recover the admission ordering

$$z_i < z_j \Rightarrow r_i \leq r_j$$

where r_i is the **oracle severity score** used in the theory.



Theoretical Analysis: Cross-Task Label Matrix

From a Single Label Vector to Stacking Across Outcomes

1) One outcome (task t)

Raw labels for task t
(binary outcome)

$$y^{(t)} = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix}$$

Admission 1
Admission 2
Admission 3
Admission 4
Admission 5
 \vdots
Admission n

Center the labels
(remove mean)

Centered label vector
 $\ell^{(t)} = Hy^{(t)}$

$$\ell^{(t)} = \begin{bmatrix} 0.48 \\ -0.52 \\ 0.48 \\ -0.52 \\ 0.48 \\ \vdots \\ -0.52 \end{bmatrix}$$

Each vertical vector $\ell^{(t)}$
is a centered label vector
for one outcome.

2) Repeat for all outcomes

T centered label vectors (one per outcome)

$$\begin{matrix} \ell^{(1)} & \ell^{(2)} & \ell^{(3)} & \ell^{(T)} \\ \begin{bmatrix} 0.48 \\ -0.52 \\ 0.48 \\ -0.52 \\ 0.48 \\ \vdots \\ -0.52 \end{bmatrix} & \begin{bmatrix} 0.41 \\ -0.59 \\ 0.41 \\ -0.59 \\ 0.41 \\ \vdots \\ -0.59 \end{bmatrix} & \begin{bmatrix} 0.55 \\ -0.45 \\ 0.55 \\ -0.45 \\ 0.55 \\ \vdots \\ -0.45 \end{bmatrix} & \begin{bmatrix} 0.30 \\ -0.70 \\ 0.30 \\ -0.70 \\ 0.30 \\ \vdots \\ -0.70 \end{bmatrix} \end{matrix}$$

n admissions (rows)

3) Stack across outcomes (matches the proof)

Stacked label matrix $\tilde{W} \in \mathbb{R}^{T \times n}$
outcomes as rows, admissions as columns

$$\tilde{W} = \begin{bmatrix} \sqrt{\beta_1} \ell^{(1)T} \\ \sqrt{\beta_2} \ell^{(2)T} \\ \vdots \\ \sqrt{\beta_T} \ell^{(T)T} \end{bmatrix} \in \mathbb{R}^{T \times n}$$

Outcome 1 (weight β_1)	0.48	-0.52	0.48	-0.52	0.48	...	-0.52
Outcome 2 (weight β_2)	0.41	-0.59	0.41	-0.59	0.41	...	-0.59
\vdots							
Outcome T (weight β_T)	0.30	-0.70	0.30	-0.70	0.30	...	-0.70

n admissions (columns)

T outcomes (rows)

$\in \mathbb{R}^{T \times n}$

Theoretical Analysis: Rank-One Projection

From W to $W^T W$ and Rank-One Signal

How stacked outcome labels induce a shared admission-level direction

1) Stacked label matrix W
(size $T \times n$)

Each row is a centered label vector for one outcome (task).

$$W = \begin{bmatrix} \sqrt{\beta_1} e^{(1)T} \\ \sqrt{\beta_2} e^{(2)T} \\ \vdots \\ \sqrt{\beta_T} e^{(T)T} \end{bmatrix}$$

	Admission 1	Admission 2	...	Admission n
Outcome 1	0.48	-0.52	...	0.12
Outcome 2	0.41	-0.59	...	-0.23
\vdots	\vdots	\vdots	...	\vdots
Outcome T	0.30	-0.70	...	0.05

n admissions

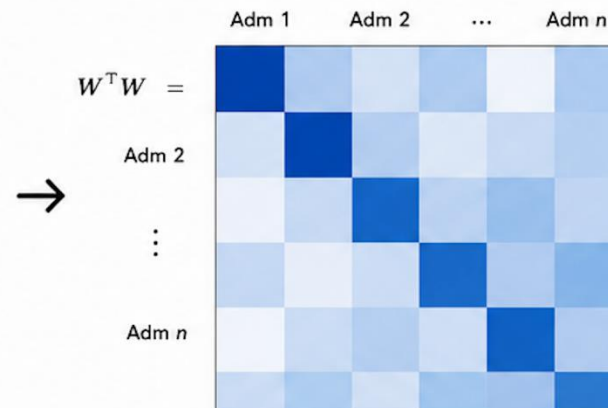
$$W \in \mathbb{R}^{T \times n}$$

T outcomes (rows) \times n admissions (columns)

2) Admission-level label-alignment matrix
 $W^T W$ (size $n \times n$)

Compares all pairs of admissions by how similarly they behave across outcomes.

$$W^T W = \sum_{t=1}^T \beta_t e^{(t)} e^{(t)T}$$



$$W^T W \in \mathbb{R}^{n \times n}$$

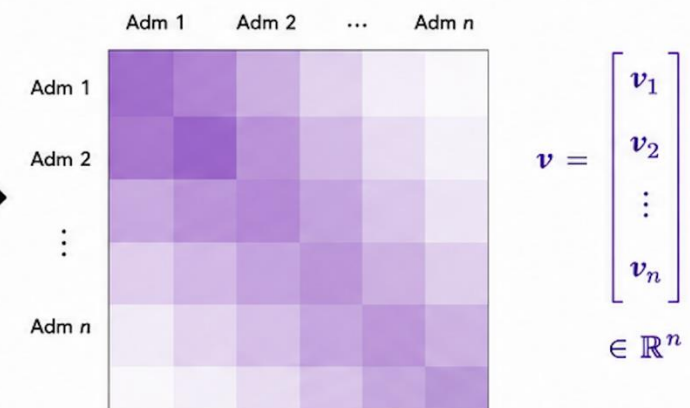
Symmetric, positive semidefinite (PSD)

3) Rank-one approximation
(shared signal)

If outcomes share a common severity axis, $W^T W$ is approximately rank one.

$$W^T W \approx \sigma_1^2 v v^T$$

(rank 1)



v is the shared admission-level direction (common severity axis across outcomes).

Linking Theory to MLCI

In the theory, we define an **oracle severity score**:

$$r_i = \text{oracle severity score for admission } i$$

It preserves the **latent severity ordering**:

$$z_i < z_j \Rightarrow r_i \leq r_j$$

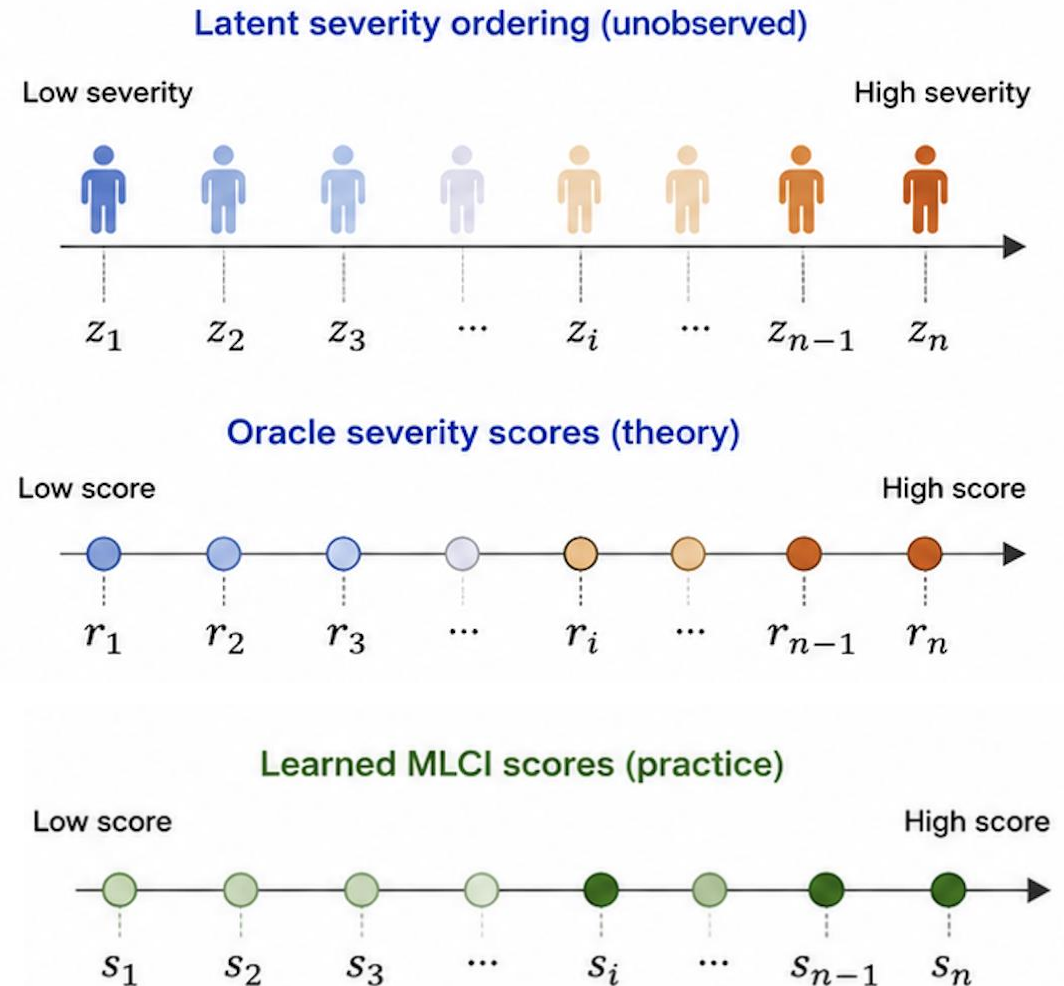
In practice, z_i and r_i are **unobserved**.

MLCI learns

$$s_i = s_\theta(X_i)$$

from diagnosis codes as a **data-driven approximation to r_i** :

$$s_i \approx r_i$$

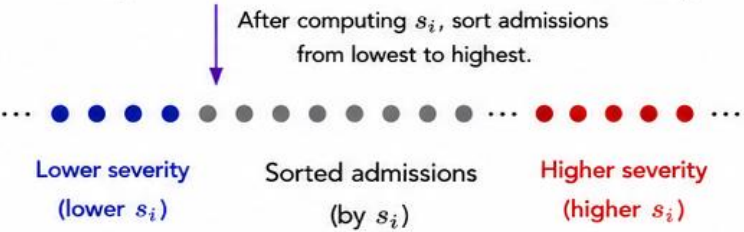
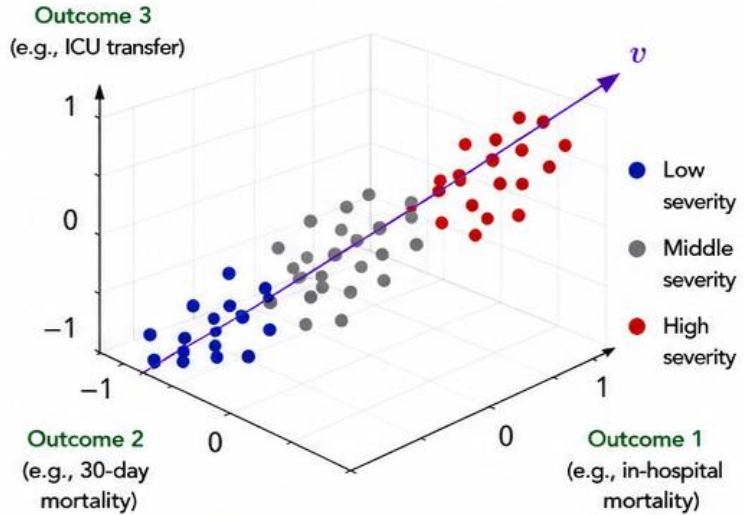


Theoretical Analysis: High Severity Threshold

The same theory also gives a **high-severity cutoff**.

1 Order admissions by the learned MLCI score s_i

Compute each admission's MLCI score $s_i = s_\theta(X_i)$ from diagnosis codes and sort from low to high.

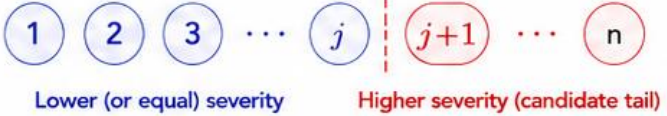


Result: a ranking of all admissions from lower to higher severity.

2 Scan possible split points

A. Setup (after sorting)

We have n admissions sorted by MLCI score:



n = number of admissions

B. What does a split point j mean?

A split point j (where $1 \leq j \leq n - 1$) divides the sorted list into two groups:

Left group (not selected): $\{1, 2, \dots, j\}$
(lower or equal severity)

Right group (candidate high-severity tail):
 $R_j = \{j + 1, \dots, n\}$

C. Scan all possible split points

We try every possible split point:

$$j = 1, 2, \dots, n - 1$$

For each j , form threshold vector g_j (1 for indices in R_j , 0 otherwise) and compute its alignment with v .

D. Choose the best split

$$j^* \in \arg \max_j \text{alignment}(v, g_j)$$

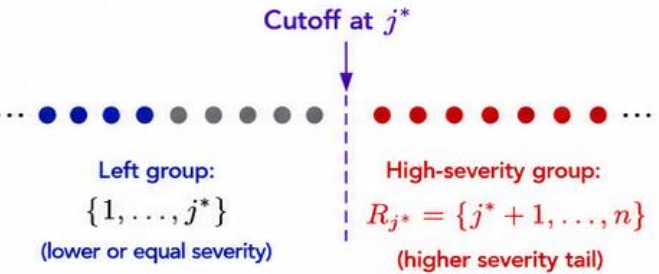
3 Choose the best split (cutoff)

The best split point is j^* .

The selected high-severity group is:

$$R_{j^*} = \{j^* + 1, \dots, n\}$$

Interpretation:
Admissions in R_{j^*} are consistently high severity across all outcomes — they align most with the shared direction v .



Threshold vector at j^* (used to evaluate alignment):

$$g_{j^*} = \left(\underbrace{0, 0, \dots, 0}_{j^* \text{ zeros}}, \underbrace{1, 1, \dots, 1}_{n - j^* \text{ ones}} \right)$$

Experimental Results: Distance Correlation (DCorr)

Table 1. Distance Correlation between each model’s unified risk score and clinical outcomes. Higher is better. Bolded values highlight the best performer per outcome. For readability, each column is scaled by the power of 10 shown in the header; divide by that factor to recover the original dCorr.

Model	MIMIC-IV				MIMIC-III			
	MORT ($\times 10^2$)	30M ($\times 10^2$)	LOS ($\times 10^2$)	ICU ($\times 10^2$)	MORT ($\times 10^2$)	30M ($\times 10^2$)	LOS ($\times 10^2$)	ICU ($\times 10^2$)
<i>Traditional Clinical Indices</i>								
Charlson (CCI)	12.59	18.44	24.55	16.09	15.54	20.26	15.28	13.07
Elixhauser (ECI)	19.98	23.87	33.15	25.98	21.38	24.70	23.10	13.38
CCI+ECI (scalar)	17.46	22.67	30.89	22.44	19.79	24.06	20.82	14.15
<i>Classical Machine Learning Baselines</i>								
kNN	22.40	23.98	30.60	34.79	22.56	23.22	23.55	11.15
Multinomial Naive Bayes	32.09	32.02	35.78	45.69	24.89	25.21	25.04	16.73
Complement Naive Bayes	32.08	32.05	35.85	45.58	24.69	25.02	23.89	16.76
FusedLogits LR	34.99	34.92	51.02	57.23	31.19	31.20	46.05	18.96
Factorization Machine (FM score)	36.41	35.82	50.80	56.77	31.96	31.28	45.38	20.51
Gradient Boosted Trees (Score/Logit)	34.54	34.58	51.00	55.61	32.28	31.88	48.16	20.17
<i>Deep Learning Baselines</i>								
Deep and Cross Network (DCN)	28.44 \pm 0.44	28.73 \pm 0.57	48.11 \pm 0.41	51.97 \pm 0.43	30.06 \pm 0.62	29.24 \pm 0.48	47.12 \pm 0.23	19.70 \pm 0.81
Deep MLP Bag-of-Codes (EmbeddingBag)	28.88 \pm 0.21	29.13 \pm 0.18	48.22 \pm 0.29	52.04 \pm 0.18	30.17 \pm 0.98	29.40 \pm 0.82	47.67 \pm 0.72	20.18 \pm 0.44
Star-GAT	26.97 \pm 1.27	27.48 \pm 1.23	47.19 \pm 0.99	50.24 \pm 0.68	29.44 \pm 0.97	29.01 \pm 0.91	48.25 \pm 0.73	20.06 \pm 0.76
Pure MIL Attention Pooling	29.25 \pm 0.62	29.53 \pm 0.43	48.37 \pm 0.36	52.20 \pm 0.63	31.11 \pm 0.70	30.40 \pm 0.31	48.11 \pm 0.43	20.41 \pm 0.53
Set Transformer (single-logit)	28.50 \pm 0.68	29.02 \pm 0.63	49.12 \pm 0.37	52.04 \pm 0.19	30.29 \pm 0.62	29.92 \pm 0.74	48.08 \pm 0.43	20.51 \pm 1.10
DeepSets (mean \oplus max pool)	28.51 \pm 0.33	28.84 \pm 0.27	48.60 \pm 0.60	51.92 \pm 0.41	29.29 \pm 0.49	29.11 \pm 0.32	48.88 \pm 0.43	21.52 \pm 0.26
Our Model	54.80 \pm 1.40	49.42 \pm 1.34	51.15 \pm 0.88	61.97 \pm 0.64	39.06 \pm 3.27	37.39 \pm 2.63	49.84 \pm 1.38	18.55 \pm 0.19

Takeaway: MLCI shows the strongest dependence for mortality outcomes and length of stay, while also achieving the best ICU-transfer dependence on MIMIC-IV.

Experimental Results: Mutual Information (MI)

Table 2. Mutual Information between each model’s unified risk score and clinical outcomes. Higher is better. Bolded values highlight the best performer per outcome. For readability, each column is scaled by the power of 10 shown in the header; divide by that factor to recover the original MI.

Model	MIMIC-IV				MIMIC-III			
	MORT ($\times 10^3$)	30M ($\times 10^3$)	LOS ($\times 10^2$)	ICU ($\times 10^2$)	MORT ($\times 10^3$)	30M ($\times 10^3$)	LOS ($\times 10^2$)	ICU ($\times 10^3$)
<i>Traditional Clinical Indices</i>								
Charlson (CCI)	9.64	18.99	3.41	1.99	16.35	26.32	0.94	14.60
Elixhauser (ECI)	20.17	28.73	6.26	3.85	26.51	30.97	3.14	14.72
CCI+ECI (scalar)	19.94	28.66	6.39	4.86	30.99	36.11	3.40	16.52
<i>Classical Machine Learning Baselines</i>								
kNN	57.70	62.71	8.43	13.06	68.28	68.07	8.55	20.65
Multinomial Naive Bayes	56.44	57.74	7.90	13.37	57.28	54.45	7.27	22.29
Complement Naive Bayes	56.68	56.99	7.75	13.26	57.75	55.43	7.12	27.53
FusedLogits LR	53.55	54.33	13.75	17.60	56.87	53.66	13.25	28.19
Factorization Machine (FM score)	54.50	54.27	13.88	16.99	62.97	57.98	13.62	27.58
Gradient Boosted Trees (Score/Logit)	64.30	63.30	14.86	17.67	65.24	66.81	14.76	24.55
<i>Deep Learning Baselines</i>								
Deep and Cross Network (DCN)	65.33 \pm 0.62	61.70 \pm 0.24	14.48 \pm 0.15	18.33 \pm 0.33	65.46 \pm 2.76	59.57 \pm 1.60	13.47 \pm 0.48	27.20 \pm 0.64
Deep MLP Bag-of-Codes (EmbeddingBag)	65.88 \pm 0.54	63.63 \pm 1.19	14.73 \pm 0.12	18.45 \pm 0.06	66.20 \pm 4.99	61.39 \pm 3.21	13.83 \pm 0.23	29.89 \pm 2.13
Star-GAT	67.38 \pm 2.48	65.32 \pm 2.87	14.89 \pm 0.27	18.89 \pm 0.47	70.31 \pm 2.27	63.74 \pm 2.44	14.10 \pm 0.54	29.28 \pm 2.68
Pure MIL Attention Pooling	67.89 \pm 0.94	65.16 \pm 2.32	14.93 \pm 0.04	18.71 \pm 0.12	71.72 \pm 1.57	64.64 \pm 0.79	14.32 \pm 0.37	26.14 \pm 2.26
Set Transformer (single-logit)	68.35 \pm 2.40	65.84 \pm 3.19	15.16 \pm 0.18	18.59 \pm 0.25	75.15 \pm 1.62	69.23 \pm 1.95	14.54 \pm 0.31	29.64 \pm 3.14
DeepSets (mean \oplus max pool)	69.92 \pm 1.26	67.22 \pm 1.60	15.36 \pm 0.11	19.24 \pm 0.19	73.69 \pm 3.29	68.08 \pm 1.88	15.17 \pm 0.30	32.42 \pm 4.22
Our Model	74.22 \pm 0.24	72.19 \pm 0.93	15.42 \pm 0.18	19.01 \pm 0.29	84.52 \pm 10.89	77.77 \pm 7.53	16.45 \pm 0.92	26.42 \pm 3.77

Takeaway: Mutual information shows the same overall trend.

Theory Diagnostics

We test the theory on MIMIC-IV and MIMIC-III using two diagnostics.

Q1: Does a shared severity signal appear in real clinical outcomes?

Rank-one alignment: Measures how much of the multi-outcome label structure is explained by one shared direction.

Table 3. Rank-one alignment and objective values on learned scores. $J(s)$ is the full multi-task nHSIC objective evaluated on the learned score vector; $J_1(s)$ is the rank-one projected objective.

Dataset	$\sigma_{1:4}$ of \widetilde{W}	$\sigma_1^2 / \ \widetilde{W}\ _F^2$	$J(s)$	$J_1(s)$
MIMIC-IV	(1.40, 1.05, 0.81, 0.53)	0.49	0.622	0.397
MIMIC-III	(1.34, 1.14, 0.84, 0.47)	0.45	0.453	0.152

Finding #1: The stacked outcome labels show moderate rank-one structure, indicating a meaningful shared severity signal across outcomes.

Q2: Can the shared direction v recover the optimal high-severity cutoff?

Threshold agreement: We compare j_v , the cutoff selected by the shared direction v , with j_J , the best cutoff for the full multi-outcome objective.

Table 4. Two-level monotone threshold diagnostics. j_J is the best split for the full multi-task threshold objective within the threshold family; j_v maximizes $\rho_j^2(v)$, equivalently selecting the best threshold for the rank-one projected objective.

Data	j_J	j_v	Tail size	$\rho_*^2(v)$
MIMIC-IV	4862	4862	137/5000 (2.74%)	0.554
MIMIC-III	11265	11265	643/11909 (5.40%)	0.220

Finding #2: The cutoff selected by the shared direction matches the best full-objective cutoff in both datasets, showing that v identifies a compact high-severity group.

Limitations of MLCI

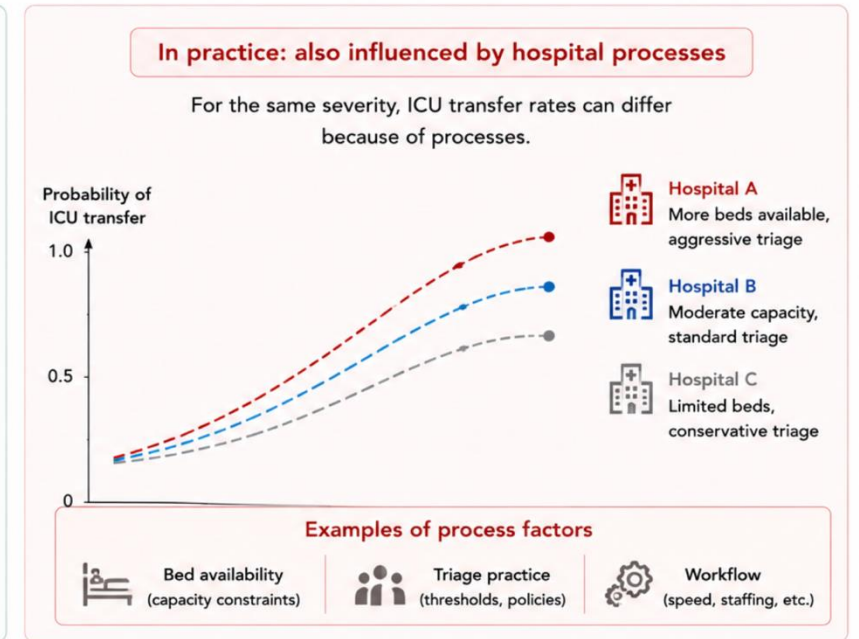
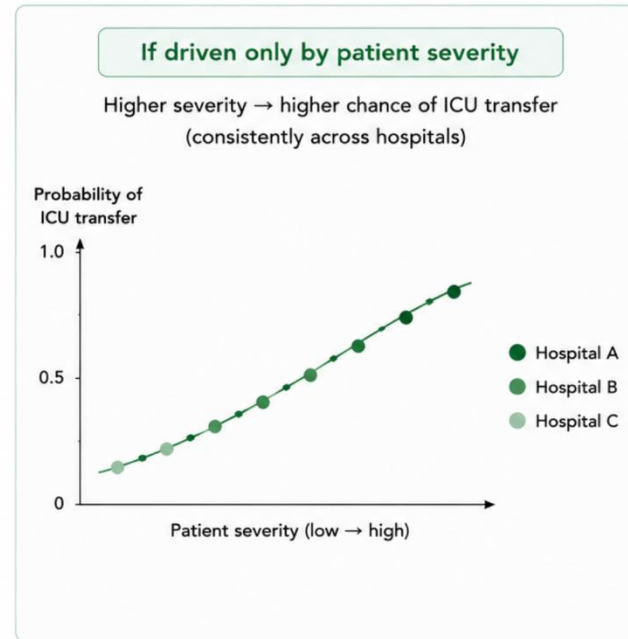
MLCI has two main limitations:

1. Shared severity assumption

MLCI assumes outcomes share a common severity signal. This assumption may be weaker for clinical outcomes driven by task-specific factors.

2. Some outcomes reflect hospital processes

For example, ICU transfer can depend on bed availability, triage practice, and workflow, rather than only patient severity.



Thank you for your time!

This project is funded by the United States National Science Foundation (NSF).



Connect
with me!

