

CAMP

Coherent Alignment of Multimodal Prototypes for Explainable Complementary Learning

A compact prototype head for complementary multimodal classification:
gate → retrieve → explain.

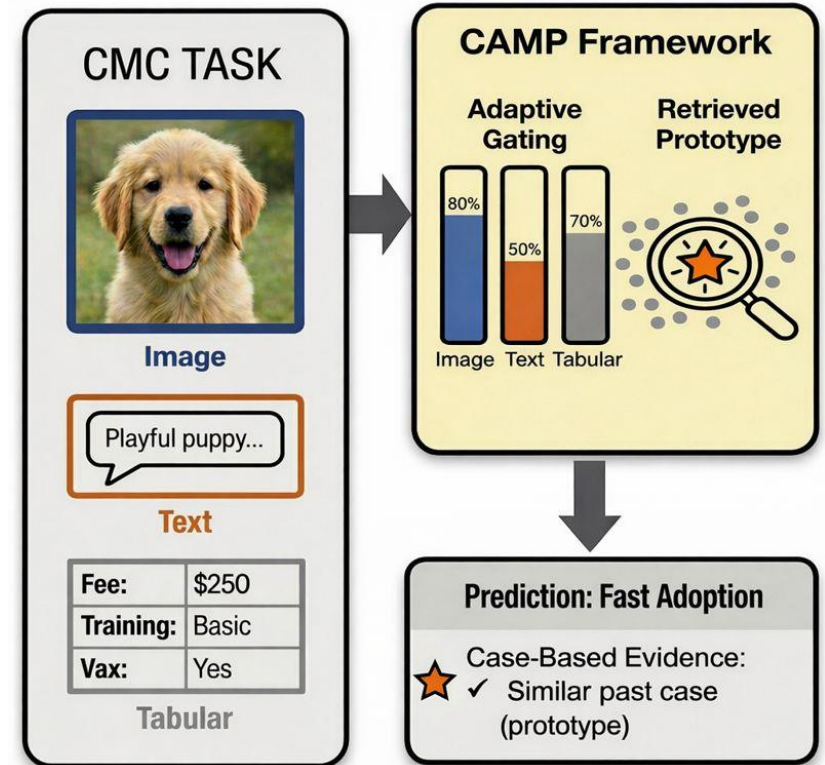
Alvaro Lopez Pellicer · Eoin M. Kenny · Simran Lamba · Shubham Sharma · Plamen P. Angelov · Saumitra Mishra

JPMorgan AI Research · Lancaster University

Disclaimer:

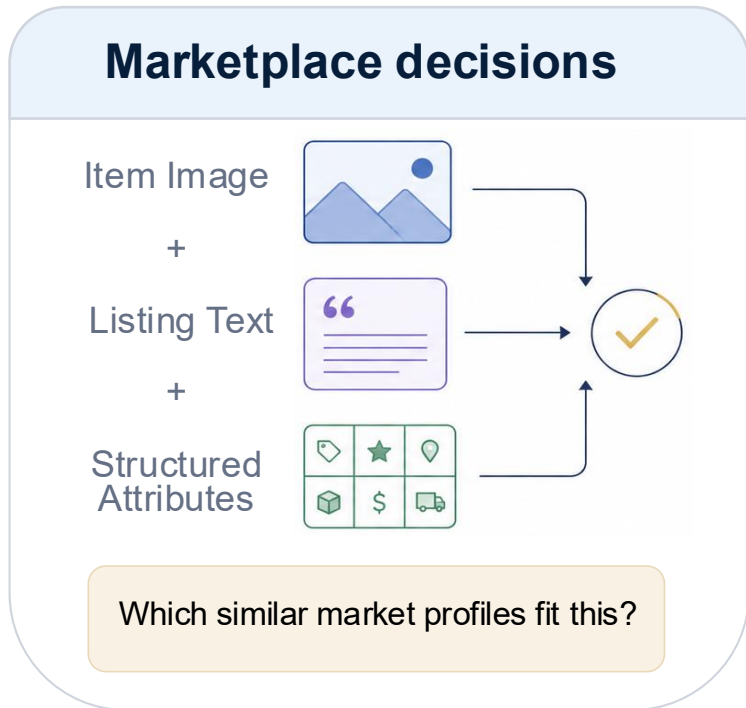
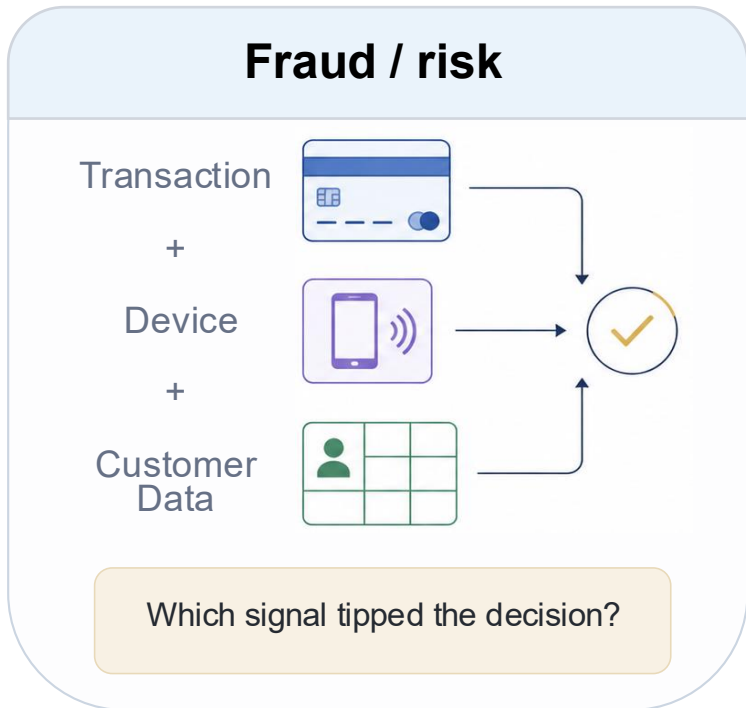
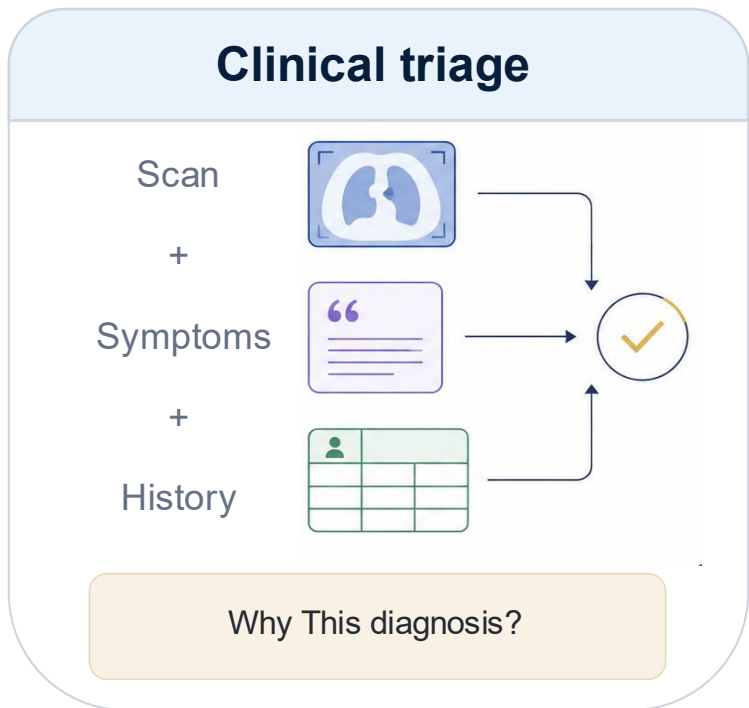
Work done during an internship at JPMorganChase AI Research.

This paper was prepared for informational purposes by the Artificial Intelligence Research group of JPMorgan Chase & Co. and its affiliates ("JP Morgan") and is not a product of the Research Department of JP Morgan. JP Morgan makes no representation and warranty whatsoever and disclaims all liability, for the completeness, accuracy or reliability of the information contained herein. This document is not intended as investment research or investment advice, or a recommendation, offer or solicitation for the purchase or sale of any security, financial instrument, financial product or service, or to be used in any way for evaluating the merits of participating in any transaction, and shall not constitute a solicitation under any jurisdiction or to any person, if such solicitation under such jurisdiction or to such person would be unlawful.



CMC is perhaps the most practical real-world multimodal setting

Important decisions often depend on evidence streams that complete one another rather than re-describe one another.



Needs: performance · auditable evidence · missing-stream robustness · case retrieval

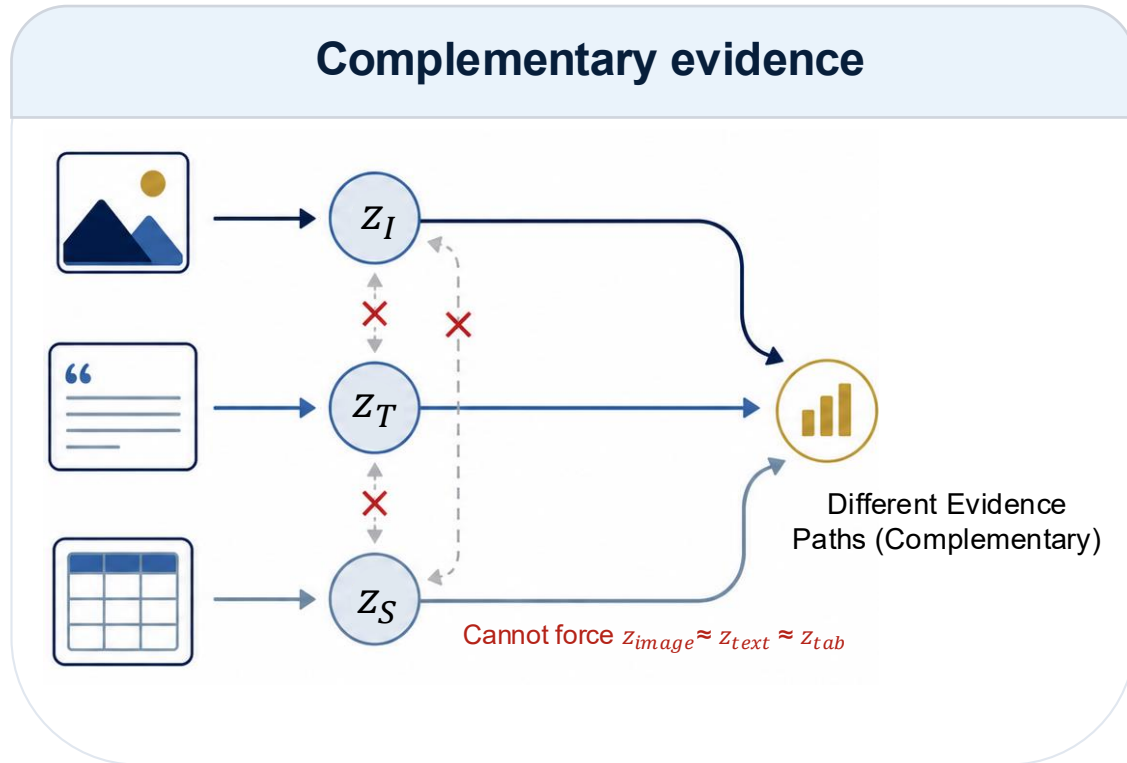
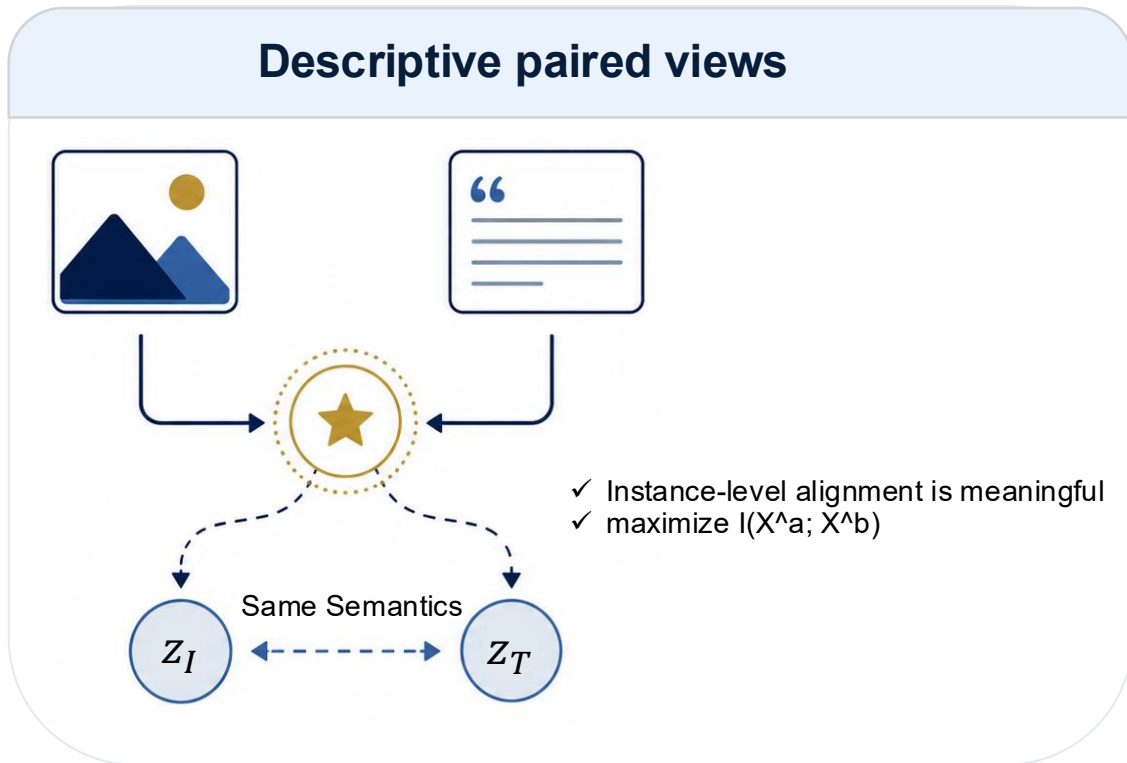
Takeaway Real CMC cases are heterogeneous evidence bundles, not paired-view captioning tasks.

Complementarity breaks paired-view alignment

In CMC, modalities can be weakly coupled but jointly predictive; forcing same-case embeddings together can erase the signal.

$$I(Y; X^{(m)} | X^{(M \setminus \{m\})}) > 0$$

At least one modality adds label information given the others

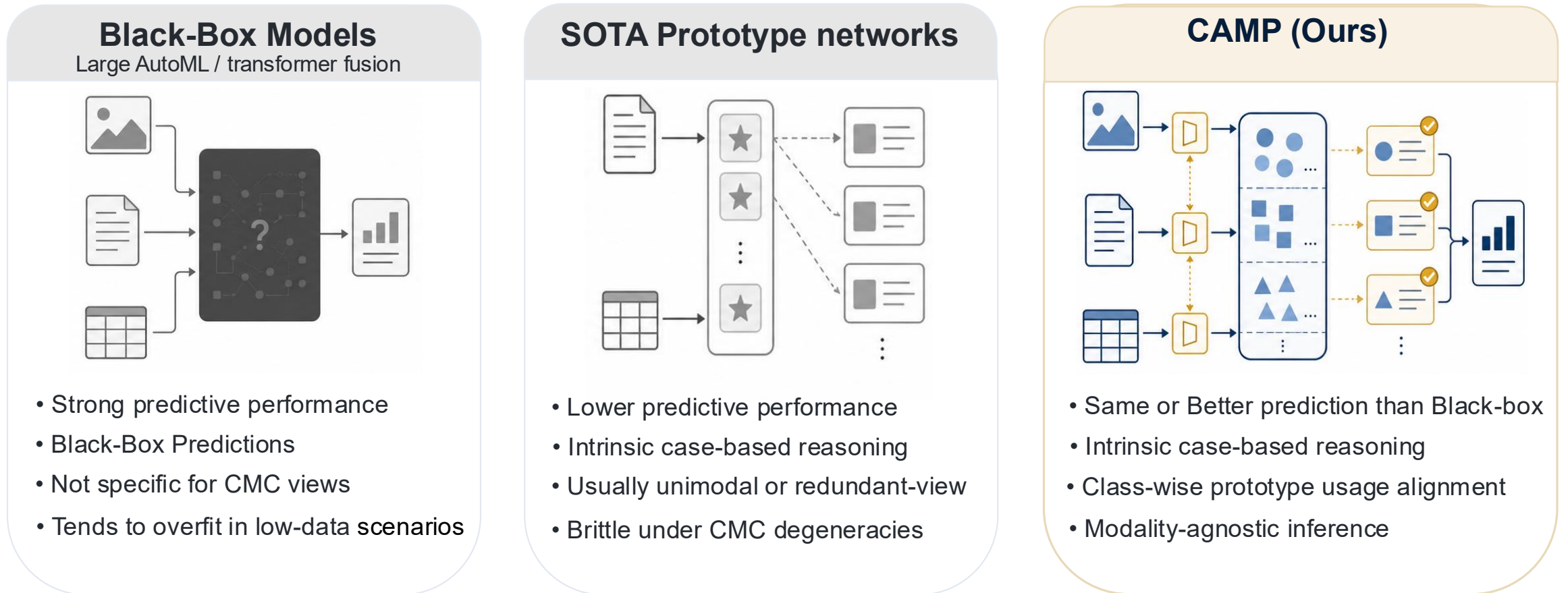


Note: z_m Refers to the embedding representation of a modality m which in this example can be I for image, T for text, or S for structured tabular records

Takeaway CMC needs coherence at the class level, not forced equality at the instance level.

The pieces exist; CMC needs the composition

Black-box fusion can be strong; prototypes can explain; CMC needs an explainable model for complementary, missing, asymmetric evidence.



Takeaway There is a clear gap for explainable-by-design CMC without sacrificing performance.

Can prototypes stay accurate, coherent, and faithful in CMC?

Yes, CAMP achieves this through gated fusion, a partitioned prototype bank, and optimal transport losses that coherently align multimodal evidence

01 Performance

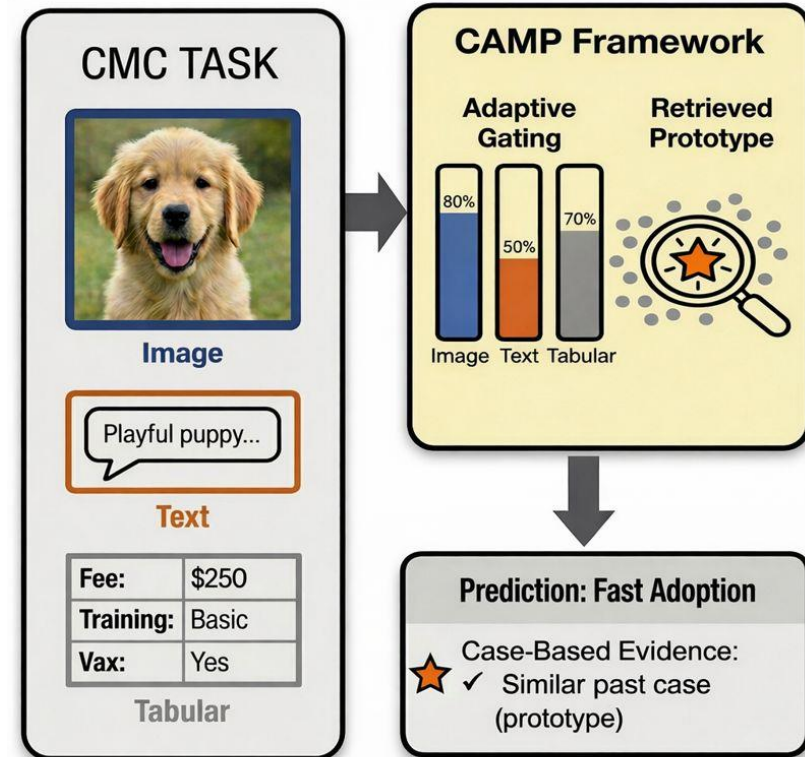
Can a compact prototype head match strong fusion baselines under fixed encoders?

02 Coherence

Can modalities stay complementary per instance while sharing class-level semantics?

03 Explainability

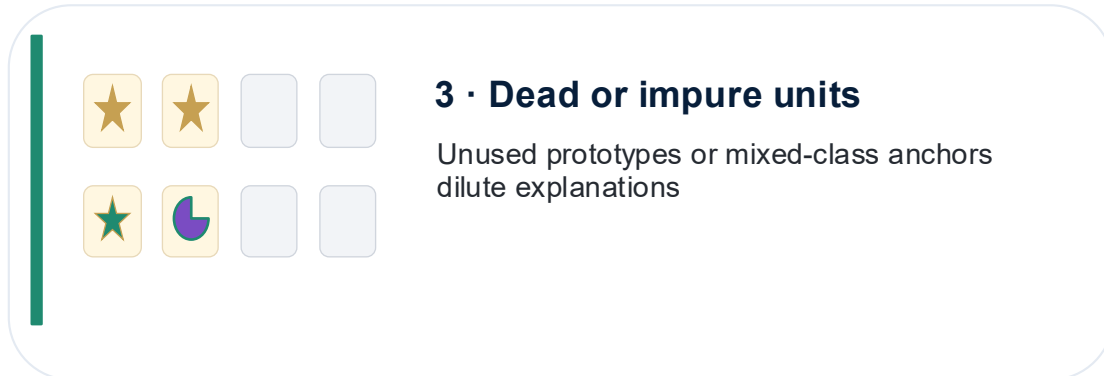
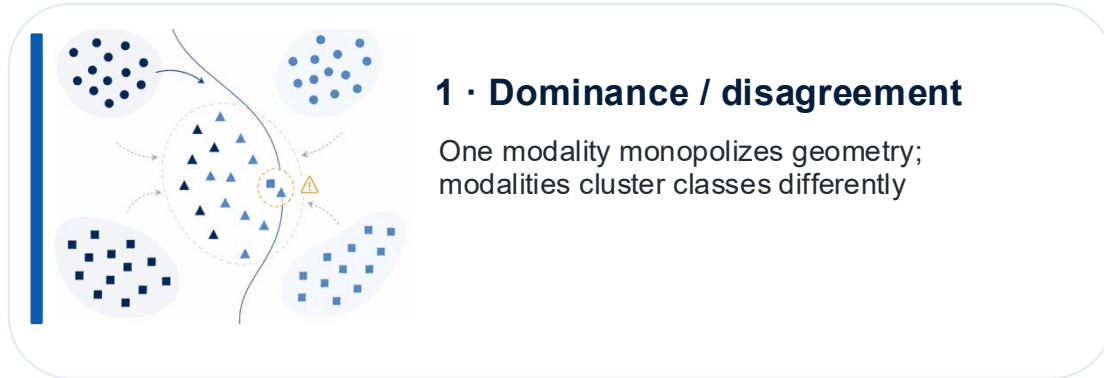
Can the same forward pass enable case-based reasoning explanations by design?



Our Method: • Align evidence usage — not individual instances • Class-Partitioned Prototype Bank • Expose gates, retrieved cases, and local evidence

Complementarity amplifies prototype failure modes

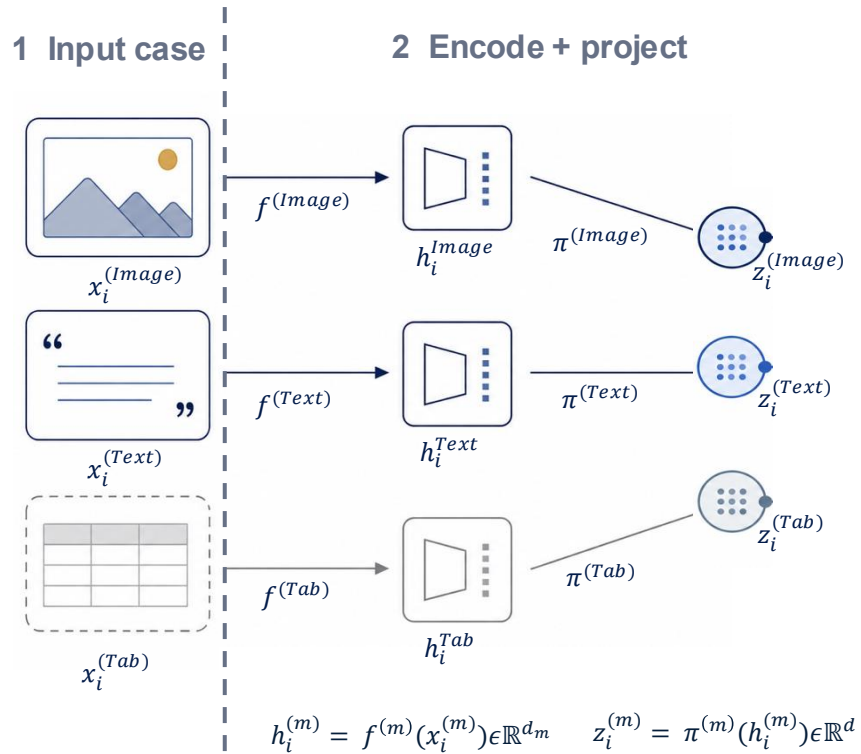
A prototype bank for CMC must be separated, used, stable, and cross-modally coherent.



Takeaway We should not merely add prototypes to a fusion model. We should make the prototype bank coherent and well-conditioned.

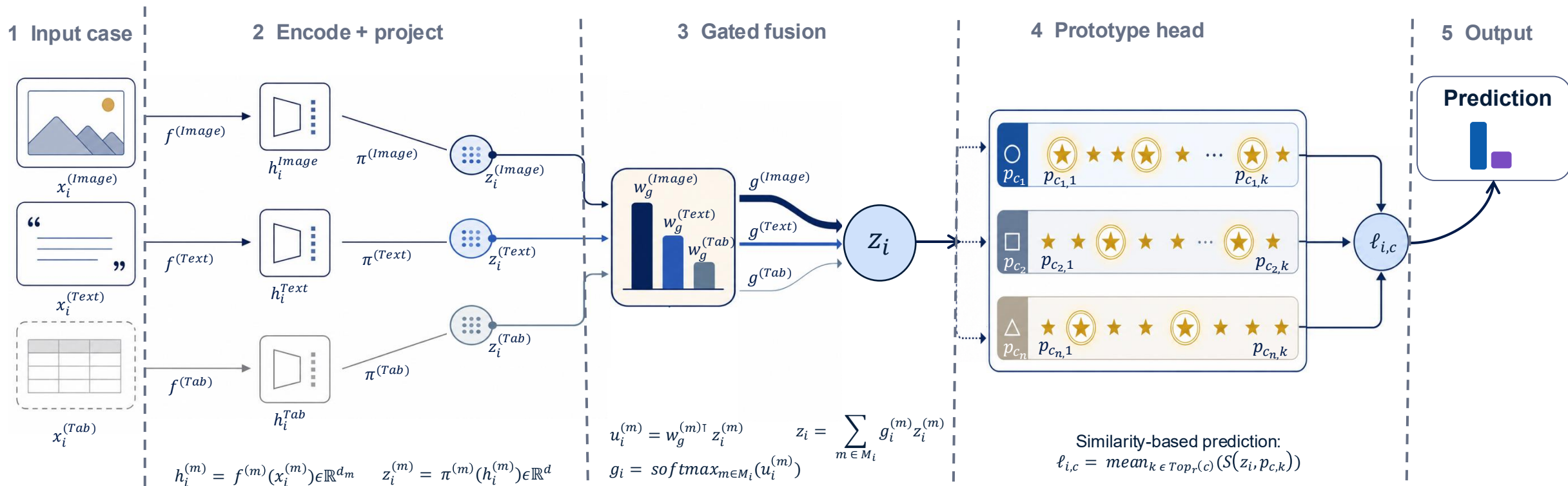
CAMP inference: encode, gate, retrieve, explain

One model handles any observed subset of image, text, and tabular modalities; explanations are returned by the same forward pass.



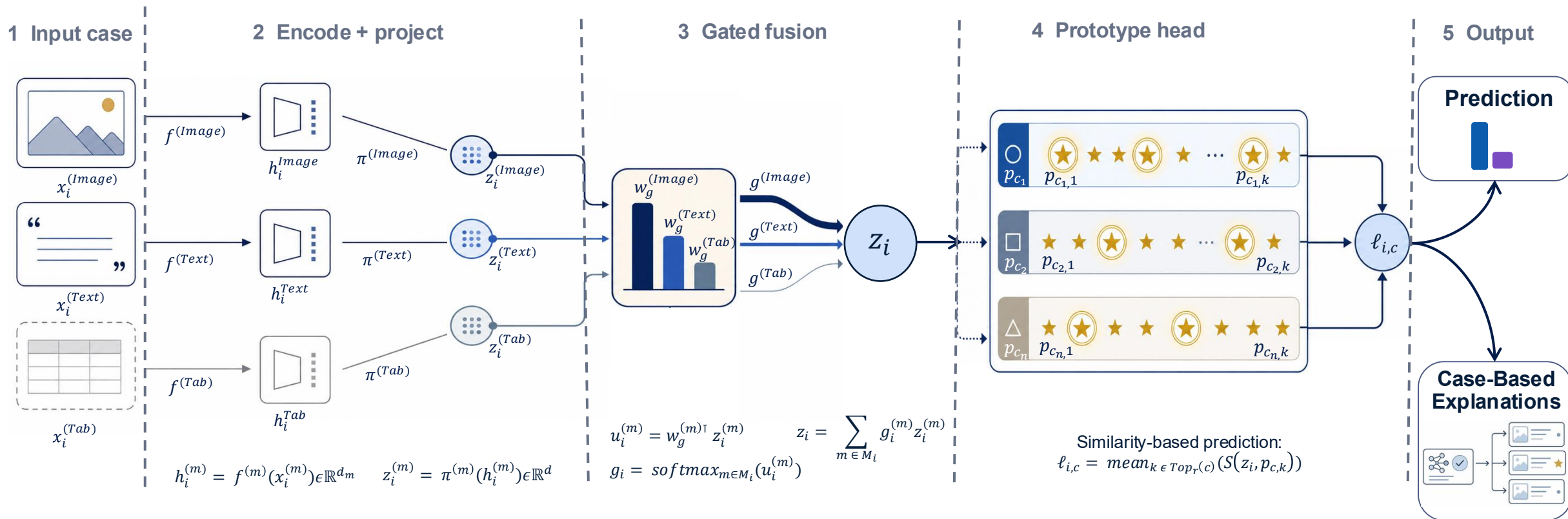
CAMP inference: encode, gate, retrieve, explain

One model handles any observed subset of image, text, and tabular modalities; explanations are returned by the same forward pass.



CAMP inference: encode, gate, retrieve, explain

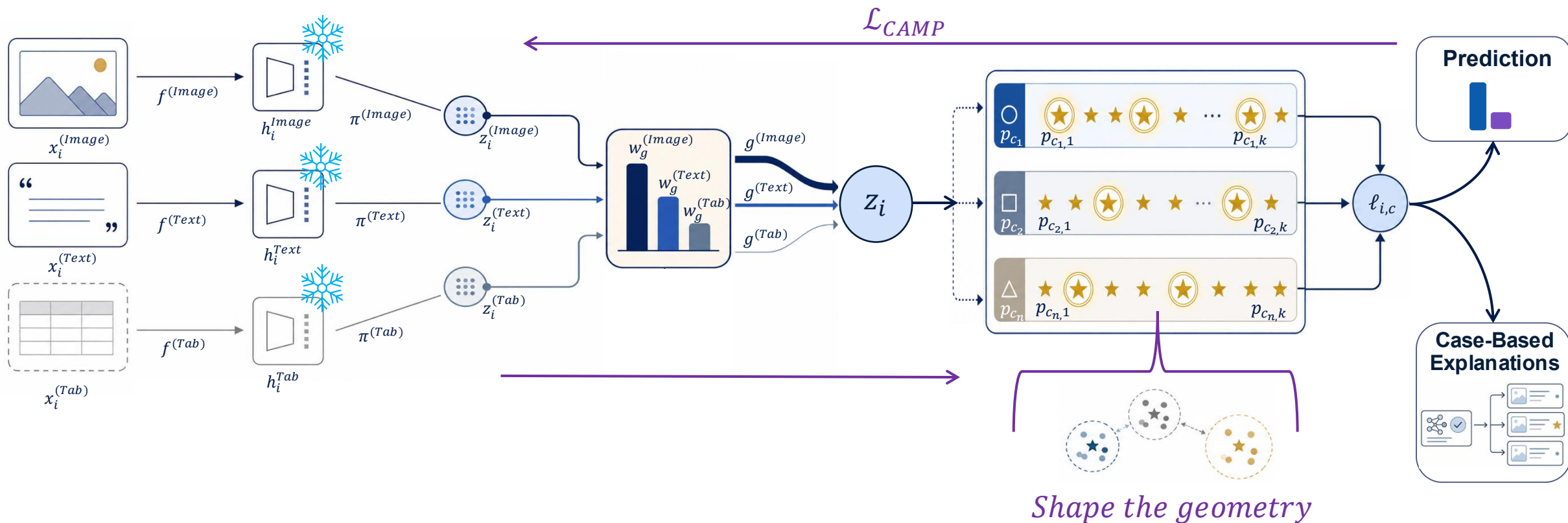
One model handles any observed subset of image, text, and tabular modalities; explanations are returned by the same forward pass.



Takeaway A single forward pass produces both prediction and case-based explanations by design.

CAMP inference: encode, gate, retrieve, explain

One model handles any observed subset of image, text, and tabular modalities; explanations are returned by the same forward pass.



Takeaway A single forward pass produces both prediction and case-based explanations by design.

Align evidence usage — not individual instances

01 Performance

Can a compact prototype head match large AutoML baselines?

02 Coherence

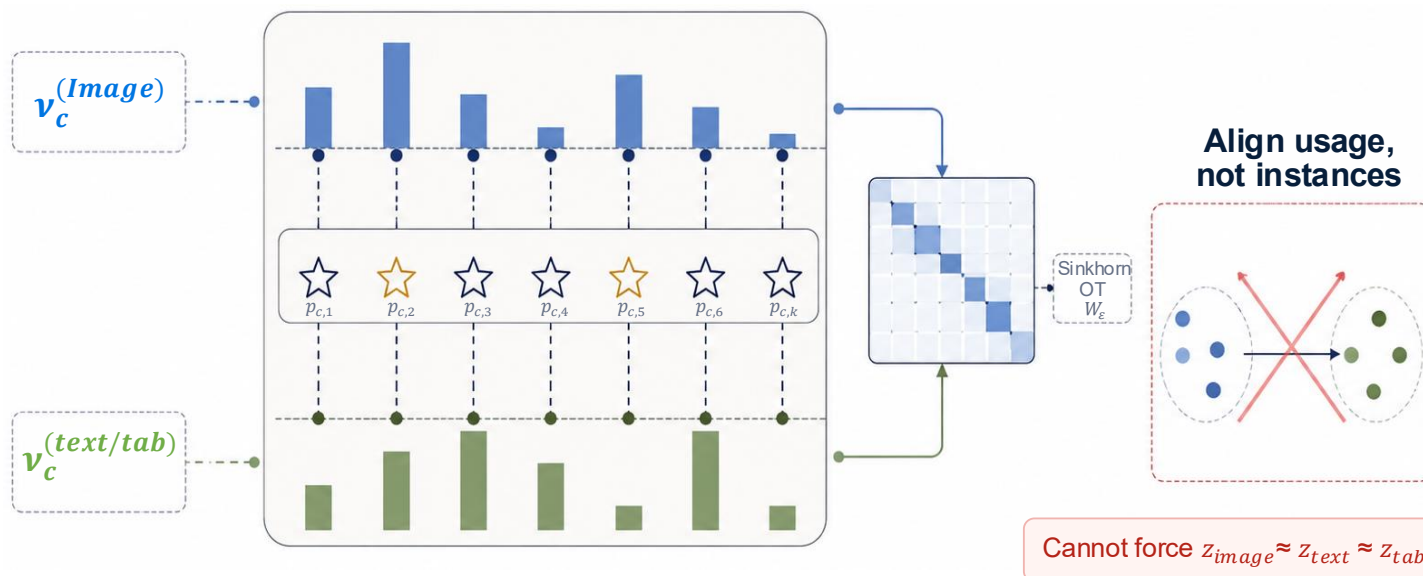
Can modalities stay complementary while sharing class semantics?

03 Explainability

Can the same forward pass enable case-based reasoning explanations by desing?

Class-wise prototype usage distributions

One class c : modalities place mass on the same prototype support



Key Takeaway Do not force modalities of the same case to coincide. Make them agree on which prototypes define a class.

Aligning Class Evidence to Match Black-Box Performance

OT synchronizes modality-specific usage distributions, supported by targeted losses that eliminate prototype degeneracies.

\mathcal{L}_{mm} : Class-wise OT

Controls: dominance / disagreement

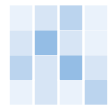
Aligns usage across modalities

For one class c , each modality induces a measure on the same prototypes

Class-wise prototype-usage measure



Multimodal class-wise OT loss



Sinkhorn OT
over prototype mass

\mathcal{L}_{geo} : margins + diversity

controls: collapse

separates classes; prevents collapse



Full training objective

$$\mathcal{L}_{CAMP} = \mathcal{L}_{pred} + \lambda_{mm}\mathcal{L}_{mm} + \lambda_{gep}\mathcal{L}_{geo} + \lambda_{alloc}\mathcal{L}_{alloc} + \lambda_{anchor}\mathcal{L}_{anchor}$$

Note CAMP's theoretical contribution is the composition: geometry, allocation, anchors, and class-wise OT turn a prototype bank into a coherent object for prediction, missing-modality inference, and case-based explanation. * Equations and guarantees are in the paper; these losses are training-time only.

Aligning Class Evidence to Match Black-Box Performance

OT synchronizes modality-specific usage distributions, supported by targeted losses that eliminate prototype degeneracies.

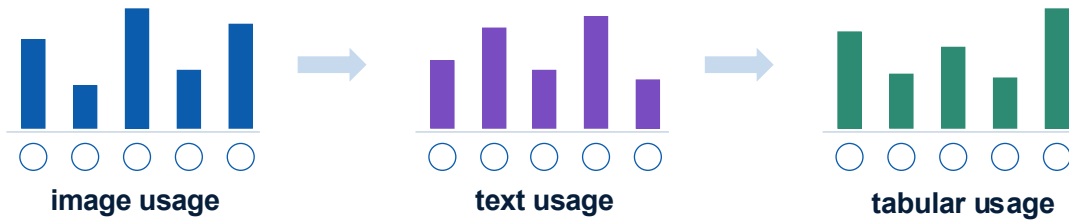
\mathcal{L}_{mm} : Class-wise OT

Aligns usage across modalities

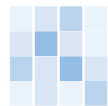
Controls: dominance / disagreement

For one class c , each modality induces a measure on the same prototypes

Class-wise prototype-usage measure



Multimodal class-wise OT loss

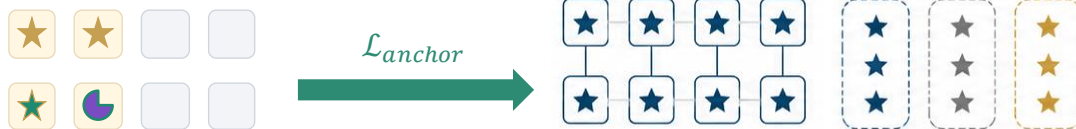


Sinkhorn OT
over prototype mass

\mathcal{L}_{alloc} : coverage + purity

controls: dead / impure units

keeps prototypes used and
class-specialized



\mathcal{L}_{geo} : margins + diversity

separates classes; prevents collapse

controls: collapse



Full training objective

$$\mathcal{L}_{CAMP} = \mathcal{L}_{pred} + \lambda_{mm}\mathcal{L}_{mm} + \lambda_{gep}\mathcal{L}_{geo} + \lambda_{alloc}\mathcal{L}_{alloc} + \lambda_{anchor}\mathcal{L}_{anchor}$$

Note CAMP's theoretical contribution is the composition: geometry, allocation, anchors, and class-wise OT turn a prototype bank into a coherent object for prediction, missing-modality inference, and case-based explanation. * Equations and guarantees are in the paper; these losses are training-time only.

Aligning Class Evidence to Match Black-Box Performance

OT synchronizes modality-specific usage distributions, supported by targeted losses that eliminate prototype degeneracies.

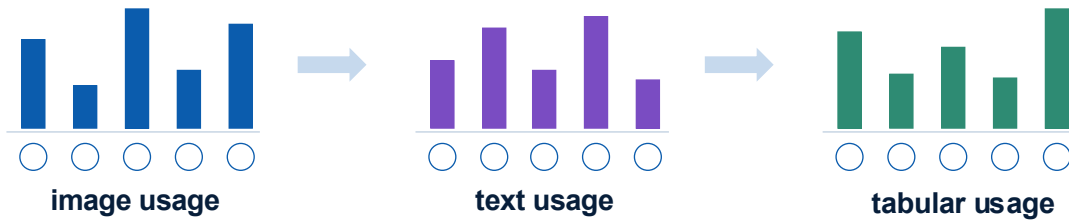
\mathcal{L}_{mm} : Class-wise OT

Aligns usage across modalities

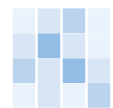
Controls: dominance / disagreement

For one class c , each modality induces a measure on the same prototypes

Class-wise prototype-usage measure



Multimodal class-wise OT loss

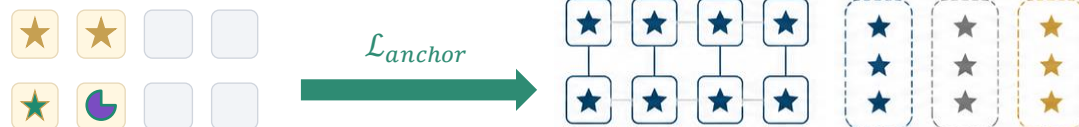


Sinkhorn OT over prototype mass

\mathcal{L}_{alloc} : coverage + purity

controls: dead / impure units

keeps prototypes used and class-specialized



\mathcal{L}_{geo} : margins + diversity

separates classes; prevents collapse

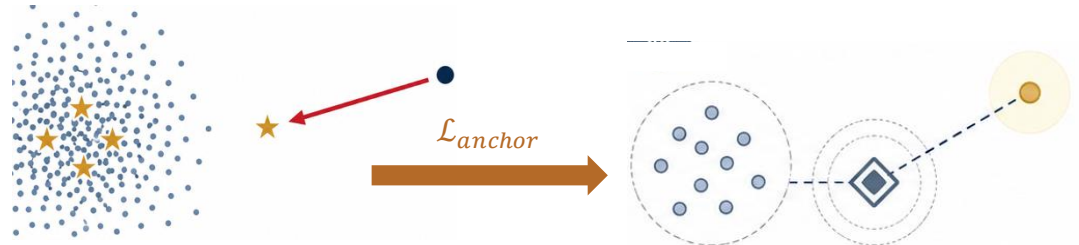
controls: collapse



\mathcal{L}_{anchor} : rare-case anchors

stabilizes outliers and drift

controls: drift



Full training objective

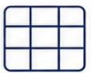




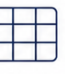


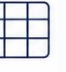
$$\mathcal{L}_{CAMP} = \mathcal{L}_{pred} + \lambda_{mm} \mathcal{L}_{mm} + \lambda_{gep} \mathcal{L}_{geo} + \lambda_{alloc} \mathcal{L}_{alloc} + \lambda_{anchor} \mathcal{L}_{anchor}$$

Note CAMP's theoretical contribution is the composition: geometry, allocation, anchors, and class-wise OT turn a prototype bank into a coherent object for prediction, missing-modality inference, and case-based explanation. * Equations and guarantees are in the paper; these losses are training-time only.

16 public CMC datasets, controlled comparisons

The benchmark separates encoder capacity from the prototype objective.

Dataset map

 	Tabular + Text 6 datasets	IMDB · FakeJob · Kick · Jigsaw · Wine · Airbnb
 	Image + Text 4 datasets	PTech · Food101 · Fakeddit · Memotion
 	Image + Tabular 3 datasets	CCD · WikiArt · HAM
  	Image + Text + Tabular 3 datasets	PetFinder · COVID · Artm

Experimental Evaluation

CAMP-F Frozen Swin-L + DeBERTa-v3 embeddings; ~660K trainable head

CAMP-E End-to-end fine-tuning after head-only warmup

Baselines XGB, LP, Proto-Raw, Proto-Base, ProtoMedX, MLP, TTT, AutoGluon

Metrics AUC/Acc/QWK; mean \pm std over five seeds; Avg and MRR

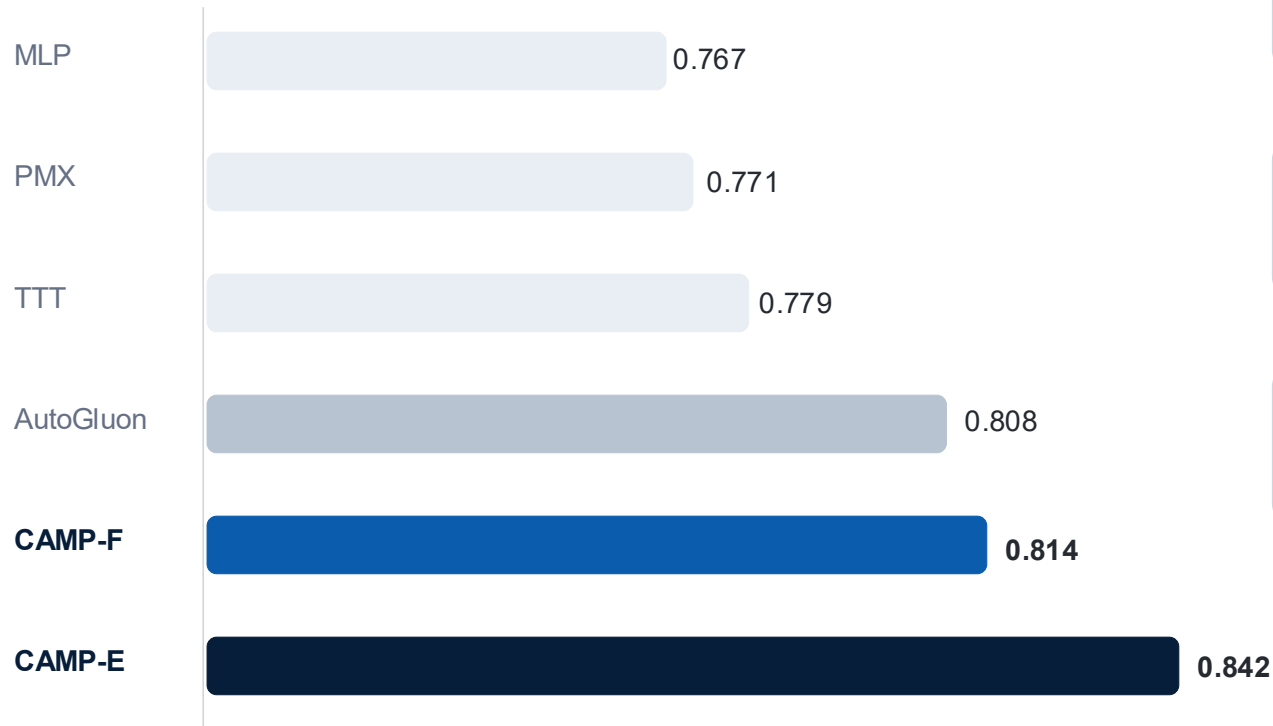
Tuning 50 Optuna trials per dataset

Note Frozen comparisons use identical embeddings to isolate the head.

Compact head, large-scale performance

CAMP-F isolates the structural objective; CAMP-E shows the ceiling with adaptation.

Suite average ↑



0.814

CAMP-F suite average

0.842

CAMP-E suite average

11/16

CAMP-F \geq AutoGluon

16/16

CAMP-E best

0.448

CAMP-F MRR vs AG
(0.296)

1.000

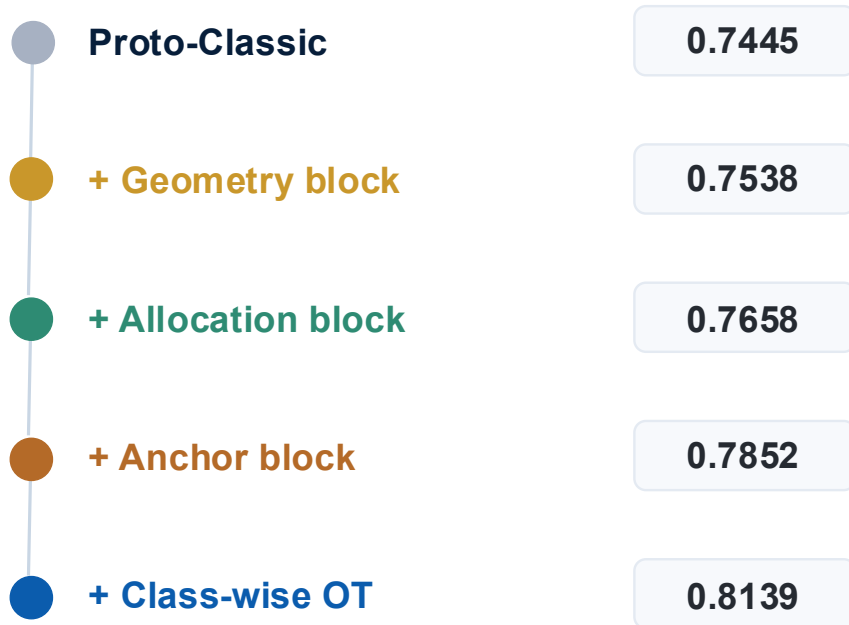
CAMP-E MRR

Key Takeaway: Frozen-encoder CAMP-E uses the same fixed embeddings as the MLP yet improves Avg 0.767 \rightarrow 0.814 and MRR 0.174 \rightarrow 0.448.

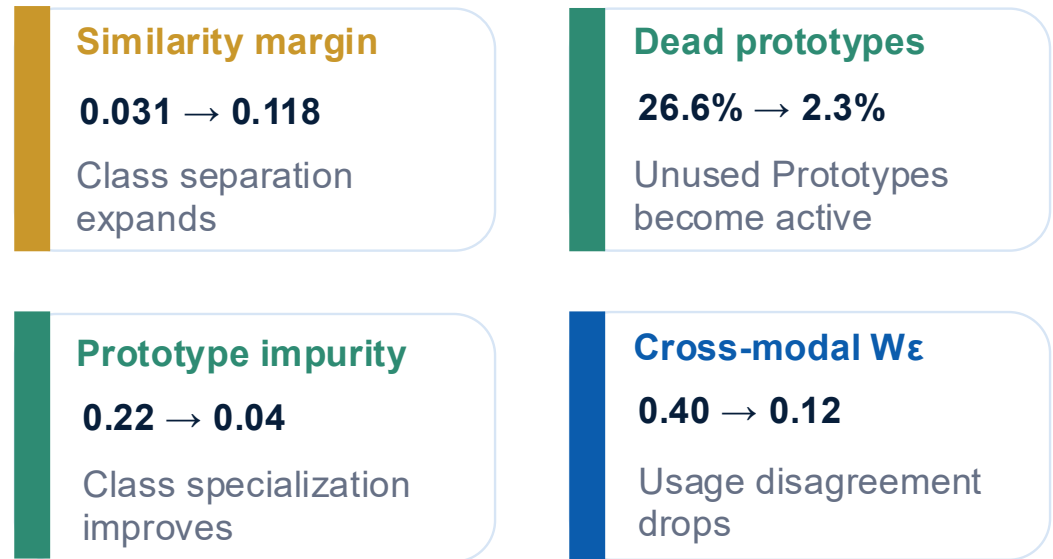
Mechanistic validation: each block fixes its failure mode

The ablation validates the theory: OT helps most after support is separated, allocated, and anchored.

Cumulative Ablations of Loss Components.



Objectives cause the intended effect:



Key Takeaway: Proto-Classic 0.7445 → CAMP full 0.8139 (+0.0694)

Largest leave-one-out hit: remove OT = -0.029

A single model degrades gracefully and explains its evidence

CAMP uses only observed modalities at inference and returns model-internal evidence.

Missing-modality robustness

No imputation or retraining

Single-drop avg **CAMP-E 0.724** AG 0.692

All configurations **CAMP-E 0.634** AG 0.610

HAM image-drop **CAMP-F 0.894** AG 0.493

PetFinder tab-drop **CAMP-F 0.425** AG 0.301

Key Takeaway: CAMP is able to work under partially observed streams.

A single model degrades gracefully and explains its evidence

CAMP uses only observed modalities at inference and returns model-internal evidence.

Missing-modality robustness

No imputation or retraining

Single-drop avg **CAMP-E 0.724** AG 0.692

All configurations **CAMP-E 0.634** AG 0.610

HAM image-drop **CAMP-F 0.894** AG 0.493

PetFinder tab-drop **CAMP-F 0.425** AG 0.301

Qualitative case study of CAMP's Multimodal explainability

Not a post-hoc surrogate: model-internal fixed-support probe reuses CAMP gates, distances, and nearest prototypes.

Airbnb case 744: Predicted Class 5 (\$100–\$120) · Confidence 98.34%

Gate Values

Observed streams only

tabular 60.84%

text 39.16%

Nearest Prototypes

★ Class 5 proto. 1

★ Class 5 proto. 2

★ Class 5 proto. 3

Fixed-support Attribution

+ num_bedrooms

+ num_bathrooms

+ "prime location"

- "Opposite of Sthybus station"

Key Takeaway:

CAMP is able to work under partially observed streams.

Attribution is performed over the model's case-based reasoning with faithfulness guarantees (See Appendix)

Interpretable CMC does not require sacrificing performance

01 Formalizes CMC

Modalities complete rather than re-describe one another

02 Coherent Alignment: Aligns usage, not instances

Class-wise OT matches prototype usage distributions

03 Explains from the model itself

Gates + retrieved cases + fixed-support evidence

04 Structural bias over scale

Sub-1M frozen head rivals >100M black-box systems; E2E improves further

Best used as a drop-in CMC head when frozen encoders, auditability, and missing-modality robustness matter.

Thank you.