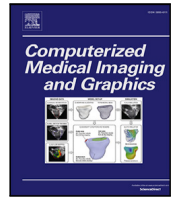




Contents lists available at ScienceDirect

# Computerized Medical Imaging and Graphics

journal homepage: [www.elsevier.com/locate/compmedimag](http://www.elsevier.com/locate/compmedimag)

## Mask-aware foundational-model embeddings for 18F-FDG-PET/CT prognosis in multiple myeloma<sup>☆</sup>

Javier Guinea-Pérez<sup>a, ID, \*</sup>, Silvia Uribe<sup>a</sup>, Sara Peluso<sup>b, c</sup>, Gastone Castellani<sup>b, c</sup>, Cristina Nanni<sup>d</sup>, Federico Álvarez<sup>a</sup>

<sup>a</sup> Universidad Politécnica de Madrid, Avenida Complutense, 30, Madrid, 28040, Madrid, Spain

<sup>b</sup> Department of Medical and Surgical Sciences, Alma Mater Studiorum-University of Bologna, Via Massarenti, 9, Bologna, Italy

<sup>c</sup> IRCCS Azienda Ospedaliero-Universitaria di Bologna, Bologna, Italy

<sup>d</sup> Nuclear Medicine, IRCCS Azienda Ospedaliero-Universitaria di Bologna, Bologna, Italy

### ARTICLE INFO

#### Keywords:

Multiple myeloma  
Survival analysis  
Radiomics  
Foundational models  
Representation learning

### ABSTRACT

**Purpose:** To test whether internal memory states from a medical foundational segmentation model can serve as compact, mask-aware embeddings for predicting progression-free survival (PFS) in multiple myeloma (MM) from whole-body [<sup>18</sup>F]FDG PET/CT, and how late fusion of PET, CT, and clinical data enhances prognostic performance.

**Methods:** We analyzed 227 newly diagnosed MM patients with PET/CT and clinical data. For two regions of interest (spine-dilated and full skeleton), we prompted MedSAM2 slice-wise using mask-derived bounding boxes and cached the final spatio-temporal memory tensor per modality. We compared two downsampling strategy to obtain per-study embeddings: channel×memory averaging with a small CNN head, and depth-attention pooling. PET and CT embeddings were combined by late fusion and passed to a DeepSurv head. We evaluated image-only and multimodal (image+clinical) models with stratified 5-fold cross-validation. The primary endpoint was Harrell's c-index (mean ± SE across folds).

**Results:** Image-only models using the averaging downsampler achieved up to 0.659±0.015 c-index (PET, spine-dilated), comparable to baseline radiomics results. Multimodal models improved discrimination to 0.710±0.032 (CT, spine-dilated), with similar performance for other PET/CT+clinical variants (0.703–0.710), improving clinical-only baselines ~ 6.5%. Averaging consistently outperformed depth-attention; concatenation and gated fusion performed comparably. PET outperformed CT within the same mask in image-only settings.

**Conclusion:** Mask-aware memory embeddings extracted from a foundational segmentation model provide effective, data-efficient imaging biomarkers for MM PFS and, when fused with routine clinical covariates, significantly improve risk stratification over clinical-only or radiomics baselines. This offers a practical path to prognostic modeling on small medical cohorts without feature design.

### 1. Introduction

Multiple myeloma (MM) is a common cancer of bone marrow (BM) plasma cells that accounts for almost 10% of all hematologic diseases, having had more than 180 000 diagnoses and 12 000 deaths in 2022 alone (Mafra et al., 2024; Ludwig et al., 2020). Few cases are diagnosed before metastasis (3% Surveillance Research Program, 2025), when disease markers are common and subtle symptoms (Rajkumar et al., 2014). Accurate risk stratification at diagnosis is central to tailoring therapy and surveillance, which can significantly improve 5-year survival rates (Surveillance Research Program, 2025).

Beyond clinical staging systems, imaging adds complementary information on tumor burden and disease biology. Whole-body [<sup>18</sup>F] fluorodeoxyglucose positron emission tomography-computed tomography (FDG PET/CT) has been included in the updated criteria for the diagnosis of MM of the International Myeloma Working Group (Rajkumar et al., 2014) as it has proven to be useful for assessing the state of the disease and treatment strategies (Agarwal et al., 2013).

Risk stratification of patients is crucial for treatment selection (Rajkumar and Kumar, 2020), but it has not yet been delegated to imaging techniques. Standardized criteria for PET/CT interpretation in MM have been proposed (Nanni et al., 2018), but visual interpretation can be

<sup>☆</sup> This article is part of a Special issue entitled: 'Data fusion in healthcare' published in Computerized Medical Imaging and Graphics.

\* Corresponding author.

E-mail address: [javier.guinea.perez@upm.es](mailto:javier.guinea.perez@upm.es) (J. Guinea-Pérez).

<https://doi.org/10.1016/j.compmedimag.2026.102752>

Received 7 November 2025; Received in revised form 27 February 2026; Accepted 6 March 2026

Available online 11 March 2026

0895-6111/© 2026 The Authors. Published by Elsevier Ltd. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

difficult and remains limited in clinical practice (Manco et al., 2023). Because of this, patient stratification is being carried out after invasive BM studies that identify chromosomal abnormalities (Kumar and Rajkumar, 2018). FDG PET has the potential to improve risk stratification by detecting medullary and extramedullary focal lesions (Terpos et al., 2011).

Recent years have seen a growing interest in the application of machine learning (ML) for FDG PET/CT interpretation in MM. Currently, images are processed using classical radiomic approaches, where handcrafted image features extract intensity, texture, and shape descriptors from segmented lesions. These features are then fed to ML algorithms with a clinical objective, be it risk stratification, biomarker discovery (Guinea-Pérez et al., 2025; Shreve et al., 2023), prognosis (Zhong et al., 2023; Ni et al., 2023), patient classification, diagnosis, and therapy response assessment (Manco et al., 2023). These studies have shown that radiomic features can correlate with tumor burden and predict outcomes, but these features rely on manually predefined descriptors that can fail to capture high order interactions or spatial context across the skeleton.

In parallel, survival analysis has taken a step forward as deep learning-based survival models have extended traditional Cox proportional-hazards to model non linear survival data (Katzman et al., 2018; Lee et al., 2018). These methods have been successfully integrated with radiomic feature extraction for survival analysis and biomarker discovery in MM and other ailments (Wiegrebe et al., 2024).

Foundational models are large-scale deep learning models trained on vast, heterogeneous datasets to acquire transferable representations that can be adapted to a wide range of downstream tasks with relatively little task-specific supervision (Bommasani et al., 2022). In computer vision, examples include CLIP (Radford et al., 2021) and SAM (Kirillov et al., 2023), which have shown that training on diverse, high-volume datasets produces general-purpose embeddings that capture both semantic and structural image properties. These embeddings can then be fine-tuned or prompted to solve specialized tasks with far less annotated data than would otherwise be required. In the medical imaging domain, foundational models are being developed by pretraining on millions of scans spanning multiple modalities (CT, MRI, PET, ultrasound) and anatomical regions (Khan et al., 2025). Adaptations of SAM, for example, have been fine-tuned to segment medical structures robustly across organ systems without retraining from scratch (Ma et al., 2024a,b).

In this work, we propose a survival analysis framework that overcomes the limitations of radiomics feature design by creating mask-aware data-efficient representations of FDG PET/CT volumes by using the internal memory state of a state of the art foundational segmentation model. These memory embeddings are generated during mask-guided volume propagation, where the model integrates anatomical prompts with imaging context. We extract and compress the final memory state into compact embeddings that capture information that radiomics descriptors cannot. We further explore late fusion strategies for PET and CT and evaluate our method in contrast to traditional radiomics features as well as measuring the improvement on a clinical feature only model.

In doing so, we address the gap between handcrafted radiomics and deep learning approaches: radiomics pipelines embed anatomical priors but are constrained by feature engineering, while deep survival networks can model complex non-linear risks but struggle to converge on small medical cohorts. Foundational medical imaging models provide a middle ground, offering robust, transferable representations that can be harnessed without retraining on small datasets. By leveraging the memory state embeddings of MedSAM2, we show that it is possible to build clinically relevant, mask-aware survival predictors from FDG PET/CT in multiple myeloma. Related work has begun to explore SAM-derived embeddings for survival in 2D medical images in multi-omics datasets (Zhu et al., 2023).

In summary, the contributions of this study are threefold: (i) we propose the use of foundational segmentation model memory states

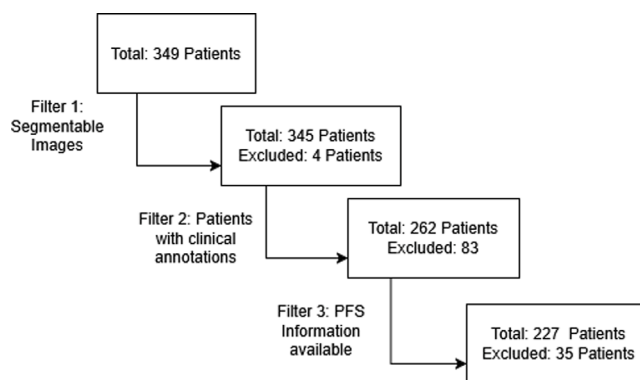


Fig. 1. Exclusion criteria of available data.

as a novel embedding for survival analysis; (ii) we benchmark PET, CT, and PET+CT late fusion embeddings against clinical and radiomics baselines; and (iii) we demonstrate that combining clinical variables with imaging embeddings yields superior progression-free survival (PFS) stratification. This positions foundational model embeddings as a promising bridge between radiomics and end-to-end deep learning in imaging-based prognostication for MM.

## 2. Materials and methods

### 2.1. Dataset

This study was based on an initial pool of 349 newly diagnosed MM patients collected since 2007 at IRCCS Azienda Ospedaliero-Universitaria Sant'Orsola-Malpighi (Bologna). All patients had baseline whole-body [ $^{18}\text{F}$ ]FDG PET/CT available. Four studies were excluded due to insufficient image quality and/or preprocessing failures that prevented reliable automatic mask generation (segmentation not obtainable), leaving 345 patients with usable imaging. From these, 262 patients had the required clinical annotations available (demographics/laboratory variables and outcome fields) and were retained for analysis; the remaining 83 were excluded due to missing clinical annotations. Finally, because PFS is the primary endpoint and therefore is not imputable, 35 patients without PFS time/event annotation were excluded, yielding a final analyzed cohort of 227 patients with PET/CT, clinical covariates, and PFS. A flow diagram summarizing these inclusion/exclusion steps and counts is provided in Fig. 1.

This dataset has, to date, not been employed in the development of any public model, and therefore all external pretrained models had no patient-level overlap with the cohort.

PFS was defined as the time from diagnosis until relapse, progression, or death, with censoring applied to patients who did not experience an event during follow-up. The available dataset does not include event subtype labels (biochemical, radiologic, or clinical progression) or an explicit adjudication protocol; therefore, we modeled PFS as recorded.

Among 227 patients, 160 (70.5%) experienced a PFS event and 67 (29.5%) were censored. The observed PFS times had a median of 30 months (25.5 months among patients with events and 35 months among censored patients). Median follow-up, estimated using the reverse Kaplan-Meier method, was 51 months. The cohort included 132 men (58.1%) and 95 women (41.9%), with a mean age of 62.5 years (standard deviation 9.7 years). Baseline staging according to the Revised International Staging System (R-ISS) (Palumbo et al., 2015) was available for most patients: 49 (21.6%) were stage I, 95 (41.9%) stage II, 25 (11.0%) stage III, while staging information was not reported for 58 (25.6%) of the cases.

The dataset included clinical covariates with established relevance to MM prognosis: age, sex, creatinine, hemoglobin, platelets, antibody

isotypes light chain ratio, and calcium. Additionally, the R-ISS stage (I/II/III) was present, which condensates information about clinical variables, some not present:  $\beta 2$  microglobulin, albumin, creatine, lactate dehydrogenase and chromosomal abnormalities. Variables with extensive missingness in this cohort (e.g., creatinine clearance) were excluded. Categorical variables were encoded and continuous variables standardized within each training fold prior to model fitting.

## 2.2. Preprocessing

PET and CT volumes were first loaded as NIfTI files and inspected. Because the two modalities differed in voxel origin and spacing despite identical qform/sform matrices, PET volumes were rigidly resampled onto the CT grid (same origin, voxel size, and dimensions) using linear interpolation with SimpleITK. No deformable registration or artifact correction was applied.

Mask generation followed the same automatic segmentation procedure described in our previous work (Guinea-Pérez et al., 2025): CT scans were segmented using the Multi-Organ Objective Segmentation (MOOSE 2.0) model (Sundar et al., 2022), which provides robust bone segmentation (median Dice coefficient 0.90 for most bones, except carpal, metacarpal and foot phalanges). From this segmentation, multiple regions of interest (ROIs) were derived through morphological operations aligned with the IMPeTUs protocol, which emphasizes bone marrow, spine, and paramedullary regions.

For the present study, we restricted the analysis to the two masks that demonstrated the best prognostic performance:

- **Spine dilated:** including vertebrae, internal spinal canal and surrounding paramedullary regions, designed to capture both intramedullary and paramedullary disease.
- **Skeleton:** a comprehensive mask encompassing the entire segmented skeleton and the dilated spine mask, designed to include total body bone disease burden.

Both masks were automatically generated on the CT images and then resampled to PET resolution, ensuring one-to-one alignment across modalities.

## 2.3. Architecture

Our pipeline converts FDG PET/CT volumes (paired with the automatically generated masks) into compact, mask-aware embeddings using the internal memory state of a foundational segmentation model, and then estimates the log-risk of progression with a DeepSurv head. PET and CT are processed in parallel with modality-specific embedding towers and are optionally combined with clinical data through late fusion before a survival head. Fig. 2 summarizes the overall architecture.

### 2.3.1. Embedding generation

We extract image representations from MedSAM2 (Ma et al., 2025), a medical adaptation of the SAM2 (Ravi et al., 2024) segmentation foundation model. The model allows for segmentation of medical images with prompts by processing them as videos. In our case, for each study, we first select a ROI mask (either the *spine dilated* or the *full skeleton*, see Section 2.2). We generate prompts by computing the tight 2D bounding box of the ROI on each axial slice containing mask pixels and propagate the segmentation top-to-bottom through the volume. To alleviate memory requirements, only one of every two mask containing slices are propagated. MedSAM2 maintains a spatio-temporal memory during this slice-wise propagation; we cache the final memory state after the last slice and use it as our representation.

Bounding-box prompts were computed per axial slice as the minimal enclosing rectangle of the CT-derived ROI mask pixels. Prompting strategy was selected to maintain stable mask propagation through the volume and reasonable segmentation performance of the downstream

MEDSAM2 model. This was verified by visual inspection on representative cases. During development we also tested alternative prompting modes (full-mask prompts and point prompts sampled within the ROI), which were less stable and/or yielded inferior downstream performance.

Concretely, for each modality (PET and CT) and ROI, MedSAM2 outputs a tensor of shape  $\mathbf{M} \in \mathbb{R}^{C \times D \times H \times W}$  (memory layers  $\times$  channels  $\times$  spatial height  $\times$  width). In order to have all embeddings have the same memory layer size, they were cropped/padded in that dimension, ending with a tensor of shape:  $\mathbb{R}^{2 \times 32 \times 64 \times 64}$ .

Given the cohort size, we consider two strategies to transform  $\mathbf{M}$  into a compact embedding:

1. **Average embedding.** We perform global averaging across the memory and channel dimensions, reducing  $\mathbf{M}$  to a  $64 \times 64$  map. This  $64 \times 64$  sized tensor was then processed through a small 2D convolutional neural network (CNN) head that maps an input with to a  $128 \times 1$  sized embedding. Starting from (1, 64, 64), the two  $2 \times 2$  max-pool layers reduce the resolution to (64, 16, 16). A final  $3 \times 3$  convolution produces (128, 16, 16), which is collapsed by global average pooling (GAP, implemented as AdaptiveAvgPool2d(1, 1)) to (128, 1, 1) and flattened to a vector in  $\mathbb{R}^{128}$ . This head is intentionally small to limit capacity and overfitting given cohort size, while GAP provides translation-robust, size-agnostic summarization of local patterns learned by the convolutions. Fig. 3(a) details the process.
2. **Light attention embedding.** Given a memory tensor  $x \in \mathbb{R}^{C \times D \times H \times W}$ , we first reduce the in-plane resolution with average 3D pooling over (1, 2, 2) to obtain  $(C, D, 32, 32)$ . A  $1 \times 1 \times 1$  pointwise 3D convolution mixes the  $C$  input channels into a compact width *base*, followed by normalization and GELU, yielding  $(base, D, 32, 32)$ .

Our depth attention block follows the squeeze-and-excitation operator (Hu et al., 2019). We squeeze spatial information by global average pooling over  $(H, W)$  to produce depth descriptors  $g \in \mathbb{R}^{base \times D}$ . A two-layer MLP with GELU then excites these descriptors to obtain slice-wise logits, and a softmax across  $D$  provides normalized attention weights  $w \in \mathbb{R}^{1 \times D}$ . The weighted sum collapses depth:

$$\tilde{x}_{c,h,w} = \sum_{d=1}^D w_d x_{c,d,h,w}, \quad \tilde{x} \in \mathbb{R}^{base \times 32 \times 32}.$$

The resulting 2D map is resized to  $S \times S$  via AdaptiveAvgPool2D, globally averaged to  $(base, 1, 1)$ , flattened to  $\mathbb{R}^{base}$ , and linearly projected to the final embedding  $\in \mathbb{R}^{out\_dim}$ . Fig. 3(b) details the attention pipeline.

*Discarded configurations.* We also evaluated three additional downsampling heads: (i) a 3D CNN head operating on  $((C \times D) \times 64 \times 64)$ , (ii) direct flattening of the  $64 \times 64$  maps, and (iii) PCA-based flattening, all under the same outer splits and training protocol. All three variants performed at chance (mean c-index  $\approx 0.50$ , high p-values) and were therefore excluded from the final pipeline.

### 2.3.2. Fusion head

We combine modality embeddings (and optionally clinical covariates) via late fusion before the survival head.

(i) *Concatenation (1D stacking).* Each modality-specific tower outputs a fixed-length vector  $\mathbf{z}_{PET}, \mathbf{z}_{CT} \in \mathbb{R}^d$ . Concatenation is performed as simple 1D stacking:

$$\mathbf{z}_{img} = [\mathbf{z}_{PET} \parallel \mathbf{z}_{CT}] \in \mathbb{R}^{2d}.$$

When including clinical variables, a vector  $\mathbf{z}_{clin} \in \mathbb{R}^{d_{clin}}$ , the final fused input to the survival head is:

$$\mathbf{z}_{fused} = [\mathbf{z}_{img} \parallel \mathbf{z}_{clin}] \in \mathbb{R}^{2d+d_{clin}}.$$

This scheme is parameter-efficient, preserves modality-specific information, and naturally supports single-modality inputs.

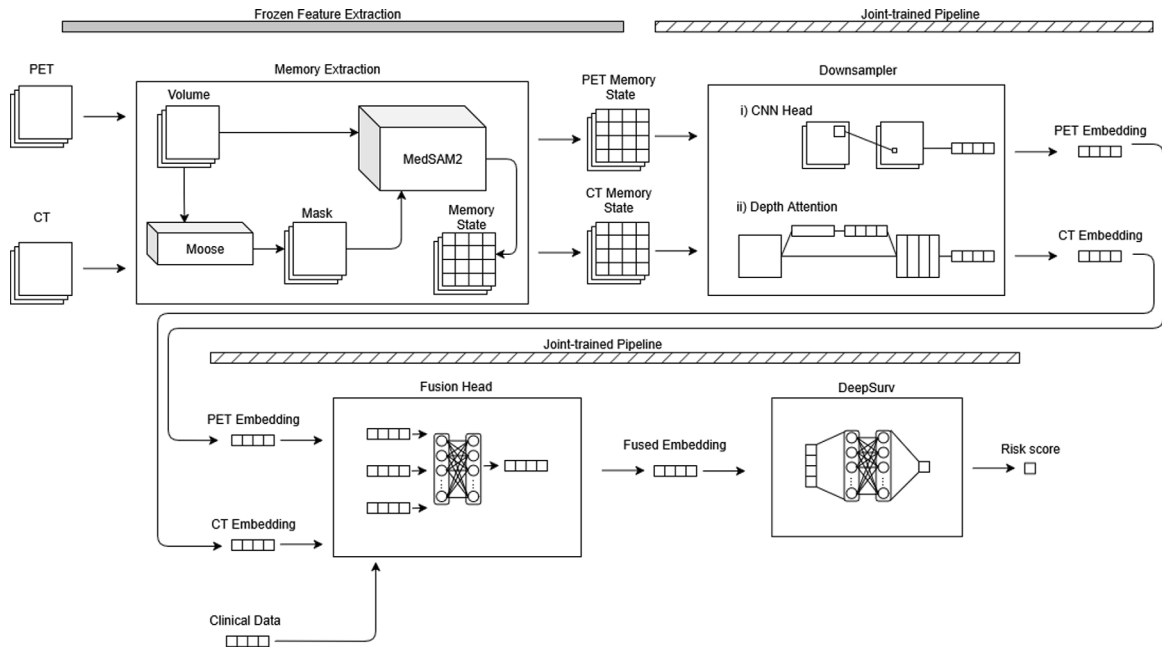


Fig. 2. Overview of the proposed pipeline. The workflow is divided into four stages: (1) Memory Extraction: For each ROI (spine dilated or skeleton) and modality (PET or CT), axial slices are prompted with mask-derived bounding boxes and propagated through MedSAM2 (Ma et al., 2025); the final memory state is cached. (2) Downsampler: The large memory tensor is compressed into a compact embedding (either with channel/memory averaging and a small CNN head and/or a lightweight depth-attention block). (3) Fusion Head: Modality-specific embeddings are late-fused (with optional clinical covariates) into a single representation either using a learned gated sum or concatenation. (4) DeepSurv: A fully connected survival head outputs the linear predictor. The top progress bar indicates trainability. Memory Extraction is frozen while Downsampler, Fusion Head and DeepSurv are trained together end-to-end.

(ii) *Lightweight gated fusion (scalar mixing)*. Alternatively, we add modality vectors  $\mathbf{z}_{\text{PET}}, \mathbf{z}_{\text{CT}} \in \mathbb{R}^d$  using a scalar gate  $\alpha \in (0, 1)$  learned via logistic regression on their 1D concatenation:

$$\alpha = \sigma(\mathbf{w}^\top [\mathbf{z}_{\text{PET}} \parallel \mathbf{z}_{\text{CT}}] + b), \quad \mathbf{z}_{\text{img}} = \alpha \mathbf{z}_{\text{PET}} + (1 - \alpha) \mathbf{z}_{\text{CT}} \in \mathbb{R}^d.$$

This gate adds only  $2d + 1$  parameters and yields an interpretable per-case modality weighting.

If clinical covariates are used, we again stack 1D:

$$\mathbf{z}_{\text{fused}} = [\mathbf{z}_{\text{img}} \parallel \mathbf{z}_{\text{clin}}] \in \mathbb{R}^{d+d_{\text{clin}}}$$

### 2.3.3. DeepSurv head

The survival head maps the fused vector  $\mathbf{z}_{\text{fused}}$  to a scalar linear predictor  $\hat{h}(x_i)$  (higher  $\hat{h}(x_i)$  = higher risk) using a deep feed-forward neural network. We update the standard DeepSurv model (Katzman et al., 2018) to include weight decay, and early stopping on validation partial likelihood. The model is trained with the average negative log partial likelihood with regularization as the loss function:

$$l(\theta) := -\frac{1}{N_{E=1}} \sum_{i: E_i=1} \left( \hat{h}_\theta(x_i) - \log \sum_{j \in \mathcal{R}(T_i)} \exp \hat{h}_\theta(x_j) \right) + \lambda \|\theta\|^2$$

where  $N_{E=1}$  is the number of observed events,  $\theta$  the parameters of the survival head,  $\mathcal{R}(T_i)$  is the risk set just prior to  $T_i$ , and  $\lambda$  controls weight decay. Dropout and batch normalization are applied to hidden layers.

## 2.4. Experiments

### 2.4.1. Experimental setup

We evaluated three families of models on PFS: (i) image-only embeddings from PET/CT, (ii) *clinical-only* baselines, and (iii) *multimodal* models combining imaging embeddings with clinical covariates. MedSAM2 memory extraction was frozen, while the downsampler, fusion, and DeepSurv head were trained end-to-end per training fold. We used stratified 5-fold cross-validation preserving event proportion and

approximate time quantiles. Optimization used AdamW cosine/step decay, dropout, BN, gradient clipping, and early stopping on validation partial likelihood. Hyperparameters were tuned by maximizing 5-fold cross-validated concordance using Sobol sequences in Optunity (Claesen et al., 2014) to search the parameter space. Because of the limited cohort size and computational cost of nested hyperparameter optimization, we used an internally tuned cross-validation procedure with a restricted search space and strong regularization. This may yield mildly optimistic performance estimates (Cawley and Talbot, 2010). Detailed search ranges and results are provided in the supplementary spreadsheet.

### 2.4.2. Models setup

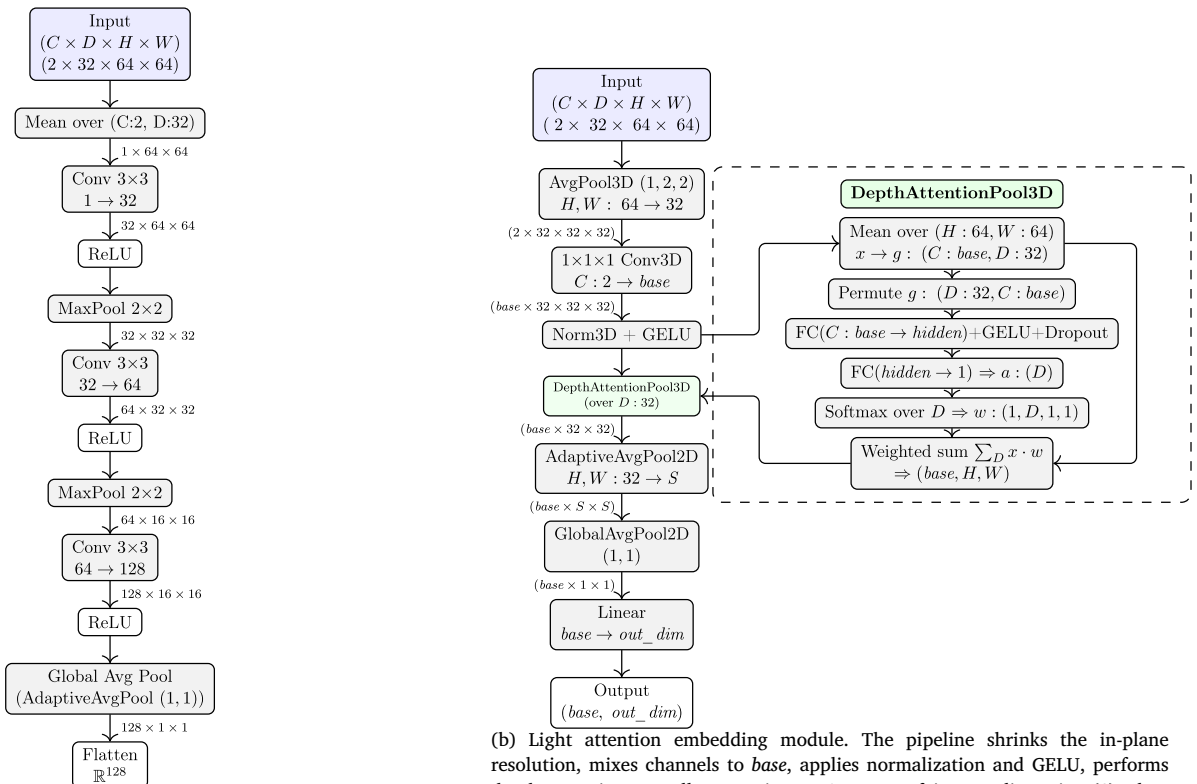
*Image-only*. We ablated (a) modality: PET, CT, PET+CT; (b) ROI: Spine-dilated vs. Skeleton; and (c) downsampler: Average embedding vs. Light attention. PET+CT used late fusion (Section 2.3.2) and additionally tested fusion strategies: concatenation vs. gated mixing. We also compare these configurations against a baseline non-medical image embedding model, ablating this way the fine-tuning and mask components of the pipeline.

*Clinical-only*. Two clinical baselines were trained using the dataset’s clinical variables (Section 2.1): CoxPH (ridge-penalized) and DeepSurv. Categorical variables were one-hot encoded, continuous variables were z-scored within each training fold and missing values were imputed from training-fold statistics.

*Multimodal*. We extend the best image-only configuration (downsampler+fusion) by adding clinical covariates. We report PET+CT+Clinical and single-modality+Clinical variants to assess complementarity with imaging. The same clinical features and preprocessing as in the clinical only model were used.

### 2.4.3. Metrics and statistical analysis

Our primary endpoint is the mean Harrell’s c-index (Harrell et al., 1982) on outer test folds, reported as mean  $\pm$  standard error (SE). For



(a) CNN embedding head used to convert the  $2 \times 32 \times 64 \times 64$  ROI map into a 128-D vector. The pipeline first applies a global average to the channel and memory dimensions, followed by *Conv-ReLU-MaxPool* twice, continued by *Conv-ReLU*, global average pooling (AdaptiveAvgPool (1, 1)), and flatten. Intermediate feature-map sizes are annotated below operations.

(b) Light attention embedding module. The pipeline shrinks the in-plane resolution, mixes channels to *base*, applies normalization and GELU, performs depth attention to collapse *D* into a 2D map of intermediate size (*S*), then applies two-stage 2D pooling and a linear projection to *out\_dim*. The inset details the depth-attention steps, where FC stands for fully connected layers.

Fig. 3. Downsampler architectures: (a) average embedding head; (b) light attention head.

the best performing model we include the Kaplan–Mayer curve of the least and most at risk patient subsets. Significance is assessed with a one-sample *t*-test against a random model (c-index of 0.5); model comparisons use paired *t*-tests on fold-wise differences (as we used identical splits).

### 3. Results

We evaluated image-only ablations, multimodal variants, and baselines. Summary metrics are in Tables 1–3.

Ablations on downsampling and fusion were conducted using image-only models (Table 1). Across origins and masks, averaging consistently outperformed attention. The top results were obtained with (i) PET+Spine Dilated+Averaging (no image fusion;  $0.659 \pm 0.015$ ,  $p = 2.08 \times 10^{-4}$ ) and (ii) Combined+Full Skeleton+Averaging with Gated fusion ( $0.659 \pm 0.018$ ,  $p = 4.84 \times 10^{-4}$ ). In contrast, Attention-based variants underperformed across all configurations (0.548–0.639 c-index).

The embedding pre-trained CNN baseline (ResNet He et al., 2016) achieved competitive performance in some image-only settings (Table 1), including one of the top configurations, but showed higher fold-to-fold variability (larger SE) than MedSAM2 memory embeddings. Across most configurations, MedSAM2 provided comparable or better discrimination with more consistent performance.

Multimodal configurations were trained using the best downsampling choice from the ablation, that is, averaging, and are summarized in Table 2. The best multimodal configuration reached  $0.710 \pm 0.032$  (CT+Spine Dilated), closely followed by Combined+Spine Dilated

( $0.710 \pm 0.034$ ) and PET+Spine Dilated ( $0.710 \pm 0.040$ ). Compared with the best image-only result, this corresponds to a relative improvement of  $\sim 7.8\%$ ; compared with the best clinical baseline (Table 3), the improvement is  $\sim 6.5\%$ .

Fig. 4 shows Kaplan–Meier curves for the first and fourth risk quartiles produced by the best overall model (clinical+CT with the spine mask). Quartiles were computed from out-of-fold risk estimates standardized using in-fold statistics. The high- vs. low-risk separation is statistically significant (log-rank  $p = 3.14 \times 10^{-3}$ ).

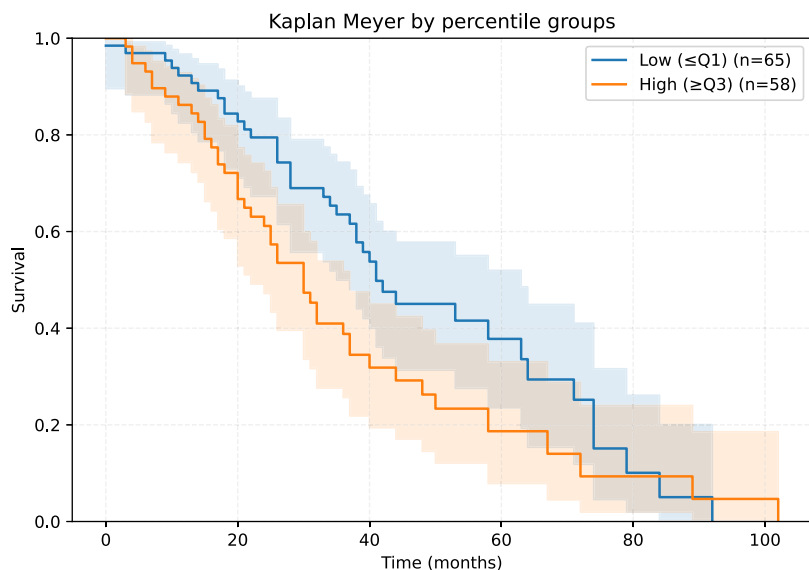
### 4. Discussion

Our novel embedding strategy proves the viability of the usage internal states of foundational medical segmentation models as information rich embeddings for survival analysis in MM. Not only are the results comparable with radiomics models (without defining the features), but results show that our embedding added to clinical models shows an improvement to clinical only models (see Table 3).

In MM, prognostic stratification can influence the overall intensity of therapy, transplant planning, and the intensity of maintenance and surveillance (Callander et al., 2025). Consequently, the risk score produced by our model is best viewed as an imaging-derived adjunct biomarker that could complement established clinical and laboratory stratification rather than replace it. Given the retrospective single-center design, limited cohort size, and incomplete availability of some guideline-grade prognostic variables in this dataset, we do not propose clinical decision thresholds or treatment changes based on our predictions; these steps require calibration and external multi-center

**Table 1**  
Image-only configurations (sorted by c-index). Reported c-indexes are mean ± SE. An en dash (–) denotes not applicable.

Origin	Mask	Model	Downsampling	Fusion	Test c-index	p-value
PET	Spine dilated	MedSAM2	Averaging	–	0.659 ± 0.015	2.081e–04
CT	Spine dilated	ResNet	–	–	0.659 ± 0.025	1.587e–03
Combined	Full skeleton	MedSAM2	Averaging	Gated	0.659 ± 0.018	4.836e–04
Combined	Full skeleton	MedSAM2	Averaging	Concatenation	0.655 ± 0.022	9.975e–04
PET	Full skeleton	MedSAM2	Averaging	–	0.654 ± 0.022	1.128e–03
CT	Spine dilated	MedSAM2	Averaging	–	0.654 ± 0.008	1.983e–05
PET	Full skeleton	ResNet	–	–	0.654 ± 0.026	1.919e–03
CT	Full skeleton	MedSAM2	Averaging	–	0.651 ± 0.021	9.451e–04
Combined	Spine dilated	MedSAM2	Averaging	Gated	0.651 ± 0.013	1.382e–04
Combined	Spine dilated	ResNet	–	Concatenation	0.650 ± 0.047	1.676e–02
Combined	Spine dilated	MedSAM2	Averaging	Concatenation	0.647 ± 0.008	3.085e–05
CT	Full skeleton	ResNet	–	–	0.641 ± 0.034	6.770e–03
Combined	Full skeleton	MedSAM2	Attention	Concatenation	0.639 ± 0.013	2.497e–04
Combined	Full skeleton	ResNet	–	Concatenation	0.638 ± 0.035	8.203e–03
Combined	Spine dilated	MedSAM2	Attention	Gated	0.638 ± 0.024	2.344e–03
Combined	Spine dilated	MedSAM2	Attention	Concatenation	0.618 ± 0.012	3.136e–04
PET	Spine dilated	ResNet	–	–	0.604 ± 0.010	2.183e–04
Combined	Full skeleton	MedSAM2	Attention	Gated	0.598 ± 0.007	9.568e–05
CT	Spine dilated	MedSAM2	Attention	–	0.586 ± 0.009	2.990e–04
CT	Full skeleton	MedSAM2	Attention	–	0.571 ± 0.015	4.956e–03
PET	Full skeleton	MedSAM2	Attention	–	0.560 ± 0.011	3.039e–03
PET	Spine dilated	MedSAM2	Attention	–	0.548 ± 0.012	8.132e–03



**Fig. 4.** Kaplan–Meier curves of the first and fourth quartile risk patients from the best performing model (clinical variables+CT with the spine mask). Risk quartiles calculated from out-of-fold risk estimations standardized with in-fold estimations. Log-rank test p-value of 3.14e–03.

**Table 2**  
Multimodal configurations (sorted by c-index).

Origin	Mask	Test c-index	p-value
CT	Spine dilated	0.710 ± 0.032	1.308e–03
Combined	Spine dilated	0.710 ± 0.034	1.697e–03
PET	Spine dilated	0.710 ± 0.040	3.104e–03
Combined	Full skeleton	0.709 ± 0.030	1.057e–03
PET	Full skeleton	0.706 ± 0.031	1.244e–03
CT	Full skeleton	0.703 ± 0.029	1.144e–03

**Table 3**  
Best models c-index (5-fold cross validation mean ± SE).

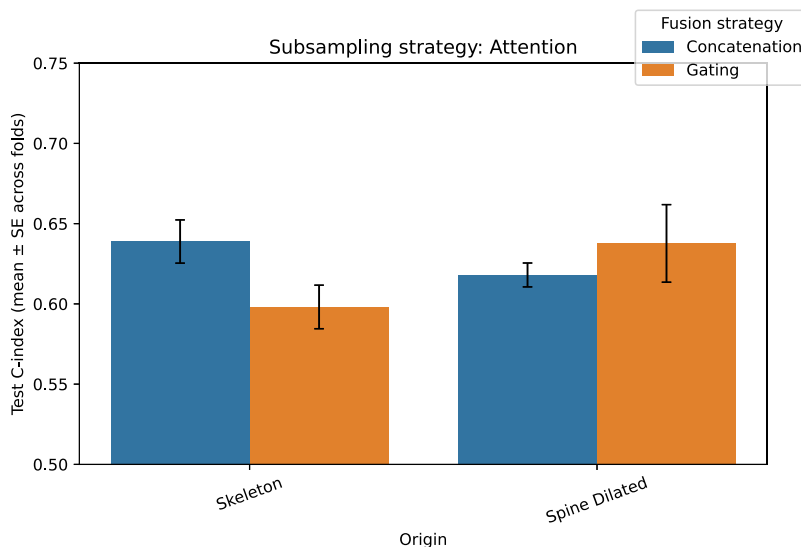
Model	c-index
Multimodal model - Our Embedding + Clinical	0.710 ± 0.030
Clinical only - DeepSurv	0.667 ± 0.035
Clinical only - Cox	0.661 ± 0.025
Image only - Our embedding	0.659 ± 0.014
Image only - Radiomics baseline	0.657 ± 0.016

validation with threshold-dependent utility analyses, as is explained in more detail in Section 5.

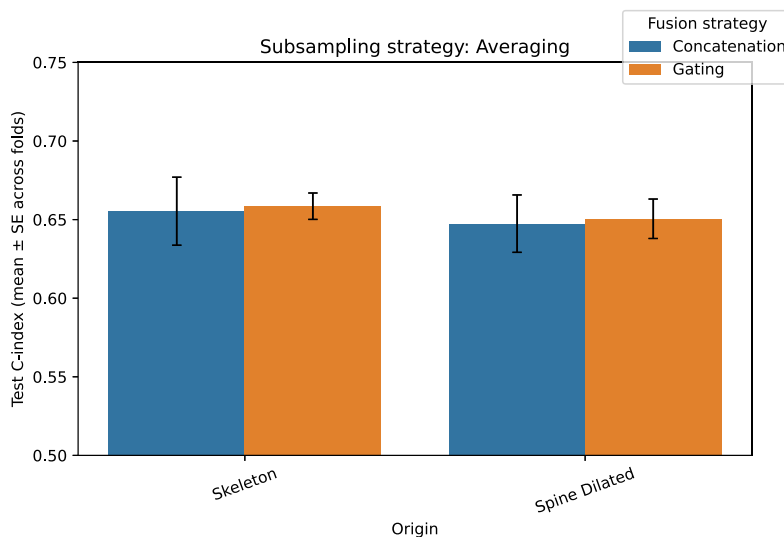
The comparison with the ResNet baseline indicates that generic visual embeddings can recover some prognostic signal from PET/CT, even without explicit segmentation prompting. However, MedSAM2 memory embeddings are explicitly conditioned on ROI prompts and integrate information accumulated during mask-guided propagation, which may

act as an inductive bias aligned with disease-relevant anatomy. In our experiments, this translated into more consistent performance across folds (lower SE) and competitive or superior c-index relative to ResNet. These findings support the view that mask-aware memory states provide a practical middle ground between handcrafted radiomics and generic pre-trained encoders when cohort sizes are limited.

In all settings (modality and mask), simple averaging across memory layers and channels outperformed the light attention downsampling



(a) Model performances by fusion strategy (hue) and mask (position) for multi-image model configurations with attention downsampling.



(b) Model performances by fusion strategy (hue) and mask (position) for multi-image model configurations with average downsampling.

Fig. 5. Comparison between both fusion strategies in imaging only models.

alternative (Fig. 6 shows the difference between downsampling strategies). Memory states are highly redundant and smooth across slices, so averaging behaves as a variance-reducing low-pass filter that suppresses slice-specific noise and propagation artifacts. Additionally, the attention head introduces additional degrees of freedom that might not be reliably constrained, increasing the risk of overfitting. Consistent with this interpretation, averaging yielded lower standard errors and p-values.

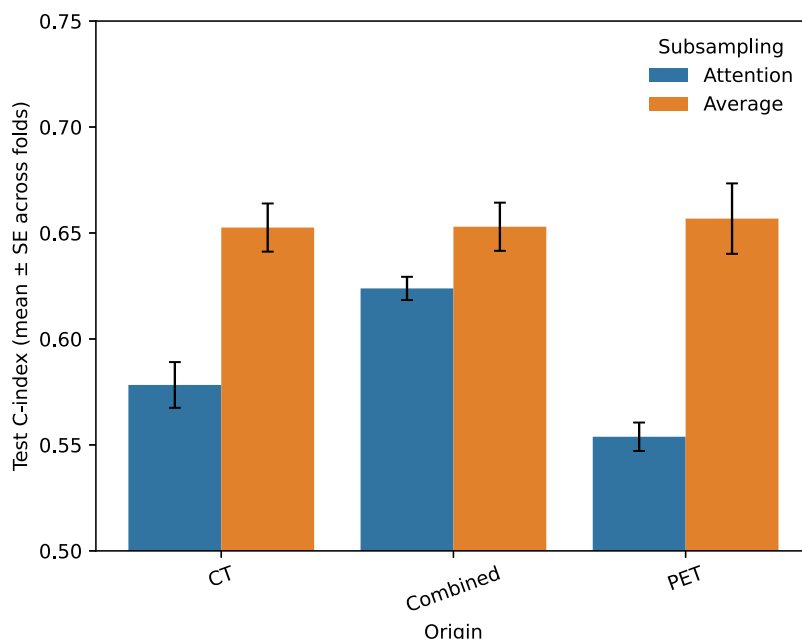
With respect to multi-image setups, fusion of PET and CT, Fig. 5 shows that 1D concatenation and scalar-gated mixing deliver statistically indistinguishable discrimination (paired tests,  $p > 0.05$ ). Both strategies access the same modality-specific vectors, and the DeepSurv head can effectively reweight concatenated features. Nevertheless, the gate remains attractive from an interpretability standpoint, as it provides a per-patient estimate of modality importance.

Within image-only models, PET consistently outperformed CT under the same mask, in line with prior work and clinical research (Guinea-Pérez et al., 2025; Zamagni et al., 2016).

Finally, mask differences were statistically insignificant both in unimodal and multimodal models (p-values  $> 0.05$ ).

### 5. Limitations and future work

Limitations of this work include a small dataset with limited clinical information. The number of events restricts statistical power, increases fold-to-fold variability, and constrains the capacity of the downsampler and fusion heads. Detailed acquisition and reconstruction metadata (scanner model, PET reconstruction settings, injected dose/uptake time, and CT parameters) were not available in this retrospective cohort, preventing explicit adjustment for scanner/protocol effects and leaving open the possibility of residual acquisition-related confounding. Additional, external, multi-center validation is needed to assess generalization under scanner/reconstruction variability. External validation would also allow for a less biased assessment, as internally tuned estimates may be mildly optimistic. Also, due to missing clinical values being imputed and because the missingness mechanism is unknown, residual bias in clinical-only and fusion models cannot be excluded.



**Fig. 6.** Aggregated model performances by downsampling strategy (hue) and origin (position) for multi-image model configurations. Reported c-index values correspond to the mean and error across folds, obtained by first averaging the c-index of each model within each fold and then aggregating these fold-wise means across the cross-validation splits.

Transferability may be limited by differences in upstream preprocessing and ROI-mask generation across sites, since our prompting is mask-driven and segmentation errors can shift the resulting memory embedding distribution. Consequently, models and hyperparameters tuned on our cohort may not be optimal under a new center’s data distribution and should be re-tuned during external validation.

Residual PET/CT misregistration (for example, due to motion or breathing) was not explicitly quantified in this study. Even though ablations show that discrimination performance is similar between unimodal sources; future work could incorporate modality alignment metrics and assess robustness to registration error.

The underperformance of the proposed lightweight attention downsampler should not be interpreted as evidence against attention mechanisms in general; higher-capacity attention architectures may behave differently but were not explored here due to sample-size constraints. With a large enough dataset, both this method or other attention mechanisms might perform better.

The retrospective dataset provides PFS as a curated time-to-event endpoint but does not include granular progression subtype labels or a documented adjudication protocol; future work using standardized event subtypes could reduce endpoint heterogeneity and enable more clinically specific modeling.

Additionally, the clinical covariates used when building the multimodal model as well as the clinical only baseline were limited. We only used available variables with acceptable completeness. Using richer multi-omics variables would likely improve performance and allow for a deeper explainability analysis.

Subgroup analyses (e.g., transplant eligibility, cytogenetic risk, extramedullary disease) were not performed because these annotations were unavailable or incomplete and subgroup sizes would be underpowered; additional robustness testing under alternative split schemes and external validation remain future work.

Our mask prompts introduce a prior over “where to look”: masks focus the model on regions deemed anatomically informative. Although prompts were derived from an automatic CT segmentation trained independently of the folds, residual concerns about prompt-induced

bias remain, as well as how much of the signal comes from the memory versus the prompt prior. Future work could compare our mask-derived prompts to coarser boxes and random prompts; and evaluate whether prompting on anatomically uninformative regions degrades performance to chance, clarifying the contribution of mask priors.

We did not perform a formal quantitative sensitivity analysis of embedding stability under prompt perturbations (bounding-box jitter/padding changes or mask boundary noise), which remains important future work to further validate robustness of the memory embeddings.

Our baseline comparisons currently include a standard pre-trained CNN (ResNet) but not transformer-based encoders (e.g., ViT) or the fine tuned image-encoder component of SAM2 (Hiera [Ryali et al., 2023](#)) in isolation. Future work should benchmark these alternatives and explicitly disentangle the contribution of the SAM2-family encoder versus prompt-guided propagation and memory updates.

The scalar-gated fusion provides per-patient modality weighting that may aid case review; projecting embedding saliency back to the masks (e.g., Grad-CAM-style analysis over the downsampler) could further expose image patterns associated with risk, a deeper explainability study could explore these ideas. In [Appendix A](#) we show an example of anatomical interpretability, though because in this pipeline there is no differentiable path between the risk prediction and the image, we have opted for exploring occlusion sensitivity.

## 6. Conclusion

This study demonstrates that mask-aware memory embeddings derived from the internal states of a foundational segmentation model can serve as effective imaging biomarkers for PFS in multiple myeloma.

Without handcrafting features, our memory-based embeddings achieved performance comparable to a strong radiomics baseline and improved a clinical-only model when fused with routine covariates ([Table 3](#)). These results support the central premise that foundational medical models encode rich, transferable information that can be repurposed—via lightweight heads and rigorous evaluation—to address prognostic tasks in small, censored cohorts.

Simple averaging across memory  $\times$  channel dimensions consistently outperforms a lightweight attention alternative, suggesting a favorable bias–variance trade-off in this data regime. PET outperforms CT under identical masks, consistent with metabolic burden driving risk, while late-fusion choices (concatenation vs. scalar gating) deliver comparable discrimination.

Overall, foundational-model memory embeddings offer a bridge between handcrafted radiomics and end-to-end deep survival networks: anatomically grounded, computationally light, and readily extensible to multimodal fusion.

### CRedit authorship contribution statement

**Javier Guinea-Pérez:** Writing – review & editing, Writing – original draft, Visualization, Software, Methodology, Investigation, Formal analysis, Conceptualization. **Silvia Uribe:** Writing – review & editing, Funding acquisition. **Sara Peluso:** Writing – review & editing, Resources, Data curation. **Gastone Castellani:** Resources, Project administration, Funding acquisition. **Cristina Nanni:** Resources, Data curation. **Federico Álvarez:** Writing – review & editing, Supervision, Project administration, Methodology, Funding acquisition, Conceptualization.

### Ethics statement

The data used in this work was gathered by IRCCS Azienda Ospedaliero-Universitaria Sant’Orsola-Malpighi in Bologna in compliance with Italian law.

Data was accessed by UPM researchers with ethics compliance following Genomed (Grant no. 101017549) and SynthIA’s (Grant No. 101172872) data sharing agreements.

### Funding

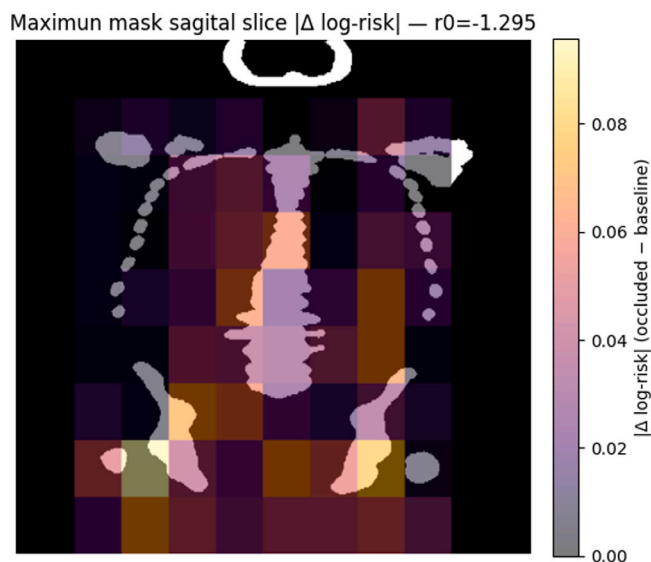
This work was supported by the European Project: SYNTHIA (Synthetic Data Generation framework for integrated validation of use cases and AI healthcare applications) Grant No. 101172872 within the Innovative Health Initiative (IHI).

### Declaration of competing interest

The authors declare the following financial interests/personal relationships which may be considered as potential competing interests: Javier Guinea-Perez reports financial support was provided by Innovative Health Initiative. Silvia Uribe reports financial support was provided by Innovative Health Initiative. Sara Peluso reports financial support was provided by Innovative Health Initiative. Gastone Castellani reports financial support was provided by Innovative Health Initiative. Cristina Nanni reports was provided by Innovative Health Initiative. Federico Alvarez reports was provided by Innovative Health Initiative. If there are other authors, they declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

### Acknowledgments

Thanks to Matteo Bastico for the idea and to Pablo Perdomo-Quinteiro, Alessandro Ceresi, Anaida Fernández García and Francisco Moreno García for their comments and support.



**Fig. A.7.** Occlusion sensitivity visualization (anatomical interpretability). Heatmap of risk sensitivity obtained by patch-wise occlusion within the prompted skeletal ROI and recomputation of the predicted log-risk. The heatmap is overlaid on the sagittal slice with the largest ROI area for a representative public whole-body PET/CT example (not used for model training). Warmer colors indicate regions where occlusion causes a larger change in predicted log-risk. In this case, the highest sensitivity is concentrated in the hip and proximal lower-limb regions within the ROI.

### Appendix A. Anatomical interpretability via occlusion sensitivity

The proposed pipeline uses MedSAM2 as a frozen embedding generator: for each study, PET/CT volumes are prompted slice-wise with mask-derived bounding boxes, MedSAM2 propagates segmentation information through the volume, and the final internal memory state is cached and compressed into a compact representation. Because the embedding extractor is not trained end-to-end and embeddings are stored offline, standard gradient-based explainability methods like GradCAM are not directly applicable to the full imaging-to-risk mapping as implemented. To obtain spatially grounded explanations under these constraints, we adopt a perturbation-based occlusion sensitivity analysis that quantifies how localized image perturbations alter the predicted risk.

Given a trained survival model (downsampler+fusion+DeepSurv) and a frozen MedSAM2 embedding extractor, we estimate an anatomical importance map by repeatedly occluding small regions of the input image, regenerating the corresponding memory embedding under the same prompting protocol, and measuring the resulting change in the model’s predicted log-risk. Regions whose occlusion causes the largest change in predicted risk are interpreted as most influential for the model’s prediction, conditional on the prompted anatomical ROI.

For each selected case, we perform the following steps:

1. **Baseline embedding and risk.** Generate the MedSAM2 memory embedding using the standard prompting protocol (slice-wise tight bounding boxes derived from the ROI mask) and compute the baseline log-risk with the trained DeepSurv head.
2. **Define an occlusion grid.** Inside the ROI, define a regular grid of patch locations. To keep the analysis anatomically constrained, we only evaluate patches whose centers lie inside the ROI mask.
3. **Occlude and re-embed.** For each patch location, modify the image by replacing voxel intensities inside the patch with a constant value (we use a neutral fill value, the within-volume mean).

We then re-run MedSAM2 propagation using the same bounding-box prompts and cache the resulting memory embedding.

- Risk delta.** Pass the occluded embedding through the same downsampler/fusion/DeepSurv head to obtain a new log-risk. The difference to the baseline log-risk (signed or absolute) is recorded as the patch importance score.
- Build an importance map.** Assign each patch score back to the occluded spatial region and aggregate overlapping patches (by averaging) to obtain a smooth voxel-level importance map. We visualize this map as overlays on representative slices or projections within the ROI.

This analysis shows which anatomical regions, when disrupted, most change the model's predicted risk. Because it is perturbation-based, it provides a model-faithful sensitivity measure for the complete pipeline. It does not imply causality, and the resulting maps depend heavily on design choices such as patch size, stride, and occlusion fill value; therefore, we interpret the visualizations qualitatively as indicators of spatial sensitivity.

Due to data-sharing restrictions for our institutional MM cohort, we cannot publish patient-level overlays from the training dataset. We therefore include an illustrative example (Fig. A.7) computed on a representative public whole-body PET/CT case (not used in model development and not diagnosed with MM) to demonstrate the procedure and verify that the resulting sensitivity patterns are spatially localized within anatomically plausible parts of the prompted ROI. The image has been sourced from The Cancer Imaging Archive (TCIA) image collection (Clark et al., 2013), and the sensitivity analysis was performed in the best performing image only model. Future work should extend this analysis to cohort-level aggregated attribution statistics under appropriate governance.

## Appendix B. Supplementary data

Spreadsheet with optimization range and best model hyperparameters available.

All code required to reproduce the proposed embedding extraction, training, and evaluation pipeline is available at <https://gitlab.com/synthiaeu/memory-embeddings>. The repository includes implementation and configuration files for MedSAM2 memory extraction, downsampling heads, late fusion, DeepSurv training, and cross-validation evaluation, occlusion sensitivity analysis, along with documentation for integrating new datasets.

Supplementary material related to this article can be found online at <https://doi.org/10.1016/j.compmedimag.2026.102752>.

## Data availability

The authors do not have permission to share data.

## References

- Agarwal, A., Chirindel, A., Shah, B.A., Subramaniam, R.M., 2013. Evolving role of FDG PET/CT in multiple myeloma imaging and management. *AJR Am. J. Roentgenol.* 200 (4), 884–890. <http://dx.doi.org/10.2214/AJR.12.9653>.
- Bommasani, R., Hudson, D.A., Adeli, E., Altman, R., Arora, S., Arx, S.v., Bernstein, M.S., Bohg, J., Bosselut, A., Brunskill, E., Brynjolfsson, E., Buch, S., Card, D., Castellon, R., Chatterji, N., Chen, A., Creel, K., Davis, J.Q., Demszky, D., Donahue, C., Doumbouya, M., Durmus, E., Ermon, S., Etchemendy, J., Ethayarajh, K., Fei-Fei, L., Finn, C., Gale, T., Gillespie, L., Goel, K., Goodman, N., Grossman, S., Guha, N., Hashimoto, T., Henderson, P., Hewitt, J., Ho, D.E., Hong, J., Hsu, K., Huang, J., Icard, T., Jain, S., Jurafsky, D., Kalluri, P., Karamcheti, S., Keeling, G., Khani, F., Khattab, O., Koh, P.W., Krass, M., Krishna, R., Kudithipudi, R., Kumar, A., Ladhak, F., Lee, M., Lee, T., Leskovec, J., Levent, I., Li, X.L., Li, X., Ma, T., Malik, A., Manning, C.D., Mirchandani, S., Mitchell, E., Munyikwa, Z., Nair, S., Narayan, A., Narayanan, D., Newman, B., Nie, A., Nibbles, J.C., Nilforoshan, H., Nyarko, J., Ogut, G., Orr, L., Papadimitriou, I., Park, J.S., Piech, C., Portelance, E., Potts, C., Raghunathan, A., Reich, R., Ren, H., Rong, F., Roohani, Y., Ruiz, C., Ryan, J., Ré, C., Sadigh, D., Sagawa, S., Santhanam, K., Shih, A., Srinivasan, K., Tamkin, A., Taori, R., Thomas, A.W., Tramer, F., Wang, R.E., Wang, W., Wu, B., Wu, J., Wu, Y., Xie, S.M., Yasunaga, M., You, J., Zaharia, M., Zhang, M., Zhang, T., Zhang, X., Zhang, Y., Zheng, L., Zhou, K., Liang, P., 2022. On the opportunities and risks of foundation models. <http://dx.doi.org/10.48550/arXiv.2108.07258>, URL <http://arxiv.org/abs/2108.07258>. arXiv:2108.07258 [cs].
- Callander, Natalie S., Mangan, Patricia A., Kumar, Shaji K., 2025. Updates to management of multiple myeloma. *J. Natl. Compr. Cancer Netw.* 23 (Supplement), <http://dx.doi.org/10.6004/jncn.2025.5008>, URL <https://jncn.org/view/journals/jncn/23/Supplement/article-e255008.xml>.
- Cawley, G.C., Talbot, N.L.C., 2010. On over-fitting in model selection and subsequent selection bias in performance evaluation.
- Claesen, M., Simm, J., Popovic, D., Moreau, Y., Moor, B.D., 2014. Easy hyperparameter search using optunity. <http://dx.doi.org/10.48550/arXiv.1412.1114>, URL <http://arxiv.org/abs/1412.1114>. arXiv:1412.1114.
- Clark, K., Vendt, B., Smith, K., Freymann, J., Kirby, J., Koppel, P., Moore, S., Phillips, S., Maffitt, D., Pringle, M., Tarbox, L., Prior, F., 2013. The cancer imaging archive (TCIA): maintaining and operating a public information repository. *J. Digit. Imaging* 26 (6), 1045–1057. <http://dx.doi.org/10.1007/s10278-013-9622-7>.
- Guinea-Pérez, J., Ceresi, A., García, A.F., Galende, B.A., Belmonte-Hernández, A., Peluso, S., Alvarez, F., 2025. Radiomics feature analysis for survival prediction in multiple myeloma: An automated PET/CT approach. *Comput. Methods Programs Biomed.* 271, 109019. <http://dx.doi.org/10.1016/j.cmpb.2025.109019>, URL <https://www.sciencedirect.com/science/article/pii/S0169260725004365>.
- Harrell, F.E., Califf, R.M., Pryor, D.B., Lee, K.L., Rosati, R.A., 1982. Evaluating the yield of medical tests. *JAMA* 247 (18), 2543–2546.
- He, K., Zhang, X., Ren, S., Sun, J., 2016. Deep residual learning for image recognition. pp. 770–778, URL [https://openaccess.thecvf.com/content\\_cvpr\\_2016/html/He\\_Deep\\_Residual\\_Learning\\_CVPR\\_2016\\_paper.html](https://openaccess.thecvf.com/content_cvpr_2016/html/He_Deep_Residual_Learning_CVPR_2016_paper.html).
- Hu, J., Shen, L., Albanie, S., Sun, G., Wu, E., 2019. Squeeze-and-excitation networks. <http://dx.doi.org/10.48550/arXiv.1709.01507>, URL <http://arxiv.org/abs/1709.01507>. arXiv:1709.01507 [cs].
- Katzman, J.L., Shaham, U., Cloninger, A., Bates, J., Jiang, T., Kluger, Y., 2018. DeepSurv: personalized treatment recommender system using a Cox proportional hazards deep neural network. *BMC Med. Res. Methodol.* 18 (1), 24. <http://dx.doi.org/10.1186/s12874-018-0482-1>.
- Khan, W., Leem, S., See, K.B., Wong, J.K., Zhang, S., Fang, R., 2025. A comprehensive survey of foundation models in medicine. <http://dx.doi.org/10.48550/arXiv.2406.10729>, URL <http://arxiv.org/abs/2406.10729>. arXiv:2406.10729 [cs].
- Kirillov, A., Mintun, E., Ravi, N., Mao, H., Rolland, C., Gustafson, L., Xiao, T., Whitehead, S., Berg, A.C., Lo, W.-Y., Dollár, P., Girshick, R., 2023. Segment anything. <http://dx.doi.org/10.48550/arXiv.2304.02643>, URL <http://arxiv.org/abs/2304.02643>. arXiv:2304.02643 [cs].
- Kumar, S.K., Rajkumar, S.V., 2018. The multiple myelomas — current concepts in cytogenetic classification and therapy. *Nat. Rev. Clin. Oncol.* 15 (7), 409–421. <http://dx.doi.org/10.1038/s41571-018-0018-y>, URL <https://www.nature.com/articles/s41571-018-0018-y>.
- Lee, C., Zame, W., Yoon, J., Schaar, M.v.d., 2018. DeepHit: A deep learning approach to survival analysis with competing risks. *Proc. AAAI Conf. Artif. Intell.* 32 (1), <http://dx.doi.org/10.1609/aaai.v32i1.11842>, number: 1. URL <https://ojs.aaai.org/index.php/AAAI/article/view/11842>.
- Ludwig, H., Novis Durie, S., Meckl, A., Hinke, A., Durie, B., 2020. Multiple myeloma incidence and mortality around the globe; interrelations between health access and quality, economic resources, and patient empowerment. *Oncol.* 25 (9), e1406–e1413. <http://dx.doi.org/10.1634/theoncologist.2020-0141>.
- Ma, J., He, Y., Li, F., Han, L., You, C., Wang, B., 2024a. Segment anything in medical images. *Nat. Commun.* 15 (1), 654. <http://dx.doi.org/10.1038/s41467-024-44824-z>, URL <https://www.nature.com/articles/s41467-024-44824-z>.
- Ma, J., Kim, S., Li, F., Baharoon, M., Asakereh, R., Lyu, H., Wang, B., 2024b. Segment anything in medical images and videos: benchmark and deployment. <http://dx.doi.org/10.48550/arXiv.2408.03322>, URL <http://arxiv.org/abs/2408.03322>. arXiv:2408.03322 [eess].
- Ma, J., Yang, Z., Kim, S., Chen, B., Baharoon, M., Fallahpour, A., Asakereh, R., Lyu, H., Wang, B., 2025. MedSAM2: Segment anything in 3D medical images and videos. <http://dx.doi.org/10.48550/arXiv.2504.03600>, URL <http://arxiv.org/abs/2504.03600>. arXiv:2504.03600 [eess].
- Mafra, A., Laversanne, M., Marcos-Gragera, R., Chaves, H.V.S., Mcshane, C., Bray, F., Znaor, A., 2024. The global multiple myeloma incidence and mortality burden in 2022 and predictions for 2045. *JNCI: J. Natl. Cancer Inst. djae321*. <http://dx.doi.org/10.1093/jnci/djae321>.
- Manco, L., Albano, D., Urso, L., Arnaboldi, M., Castellani, M., Florimonte, L., Guidi, G., Turra, A., Castello, A., Panareo, S., 2023. Positron emission tomography-derived radiomics and artificial intelligence in multiple myeloma: State-of-the-art. *J. Clin. Med.* 12 (24), 7669. <http://dx.doi.org/10.3390/jcm12247669>, URL <https://www.mdpi.com/2077-0383/12/24/7669>.
- Nanni, C., Versari, A., Chauvie, S., Bertone, E., Bianchi, A., Rensi, M., Bellò, M., Gallamini, A., Patriarca, F., Gay, F., Gamberi, B., Ghedini, P., Cavo, M., Fanti, S., Zamagni, E., 2018. Interpretation criteria for FDG PET/CT in multiple myeloma (IMPeTUS): final results. *IMPeTUS (Italian myeloma criteria for PET USE)*. *Eur. J. Nucl. Med. Mol. Imaging* 45 (5), 712–719. <http://dx.doi.org/10.1007/s00259-017-3909-8>.

- Ni, B., Huang, G., Huang, H., Wang, T., Han, X., Shen, L., Chen, Y., Hou, J., 2023. Machine learning model based on optimized radiomics feature from 18F-FDG-PET/CT and clinical characteristics predicts prognosis of multiple myeloma: A preliminary study. *J. Clin. Med.* 12 (6), 2280. <http://dx.doi.org/10.3390/jcm12062280>.
- Palumbo, A., Avet-Loiseau, H., Oliva, S., Lokhorst, H.M., Goldschmidt, H., Rosinol, L., Richardson, P., Caltagirone, S., Lahuerta, J.J., Facon, T., Bringhen, S., Gay, F., Attal, M., Passera, R., Spencer, A., Offidani, M., Kumar, S., Musto, P., Lonial, S., Petrucci, M.T., Orłowski, R.Z., Zamagni, E., Morgan, G., Dimopoulos, M.A., Durie, B.G.M., Anderson, K.C., Sonneveld, P., San Miguel, J., Cavo, M., Rajkumar, S.V., Moreau, P., 2015. Revised international staging system for multiple myeloma: a report from international myeloma working group. *J. Clin. Oncol.: Off. J. Am. Soc. Clin. Oncol.* 33 (26), 2863–2869. <http://dx.doi.org/10.1200/JCO.2015.61.2267>.
- Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., Krueger, G., Sutskever, I., 2021. Learning transferable visual models from natural language supervision. <http://dx.doi.org/10.48550/arXiv.2103.00020>, URL <http://arxiv.org/abs/2103.00020>. arXiv:2103.00020 [cs].
- Rajkumar, S.V., Dimopoulos, M.A., Palumbo, A., Blade, J., Merlini, G., Mateos, M.-V., Kumar, S., Hillengass, J., Kastritis, E., Richardson, P., Landgren, O., Paiva, B., Dispenzieri, A., Weiss, B., LeLau, X., Zweegman, S., Lonial, S., Rosinol, L., Zamagni, E., Jagannath, S., Sezer, O., Kristinsson, S.Y., Caers, J., Usmani, S.Z., Lahuerta, J.J., Johnsen, H.E., Beksac, M., Cavo, M., Goldschmidt, H., Terpos, E., Kyle, R.A., Anderson, K.C., Durie, B.G.M., Miguel, J.F.S., 2014. International Myeloma Working Group updated criteria for the diagnosis of multiple myeloma. *Lancet Oncol.* 15 (12), e538–e548. [http://dx.doi.org/10.1016/S1470-2045\(14\)70442-5](http://dx.doi.org/10.1016/S1470-2045(14)70442-5), URL [https://www.thelancet.com/article/S1470-2045\(14\)70442-5/fulltext](https://www.thelancet.com/article/S1470-2045(14)70442-5/fulltext).
- Rajkumar, S.V., Kumar, S., 2020. Multiple myeloma current treatment algorithms. *Blood Cancer J.* 10 (9), 94. <http://dx.doi.org/10.1038/s41408-020-00359-2>, URL <https://www.nature.com/articles/s41408-020-00359-2>.
- Ravi, N., Gabeur, V., Hu, Y.-T., Hu, R., Ryali, C., Ma, T., Khedr, H., Rädle, R., Rolland, C., Gustafson, L., Mintun, E., Pan, J., Alwala, K.V., Carion, N., Wu, C.-Y., Girshick, R., Dollár, P., Feichtenhofer, C., 2024. SAM 2: Segment anything in images and videos. <http://dx.doi.org/10.48550/arXiv.2408.00714>, URL <http://arxiv.org/abs/2408.00714>. arXiv:2408.00714 [cs].
- Ryali, C., Hu, Y.-T., Bolya, D., Wei, C., Fan, H., Huang, P.-Y., Aggarwal, V., Chowdhury, A., Poursaeed, O., Hoffman, J., Malik, J., Li, Y., Feichtenhofer, C., 2023. HierA: a hierarchical vision transformer without the bells-and-whistles. In: *Proceedings of the 40th International Conference on Machine Learning*. In: ICML'23, vol. 202, JMLR. org, Honolulu, Hawaii, USA, pp. 29441–29454.
- Shreve, J., Charalampous, C., Warsame, R., Pritchett, J., Kapoor, P., Buadi, F., Paludo, J., Dingli, D., Shah, M., Fonder, A., Hayman, S., Leung, N., Gertz, M., Kyle, R., Broski, S., Rajkumar, V., Kumar, S., Binder, M., 2023. P921: Artificial intelligence based FDG PET/CT radiomics for risk stratification in newly diagnosed multiple myeloma. *HemaSphere* 7 (Suppl), e156011e. <http://dx.doi.org/10.1097/01.HS9.0000970588.15601.1e>, URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC10431033/>.
- Sundar, L.K.S., Yu, J., Muzik, O., Kulterer, O.C., Fueger, B., Kifjak, D., Nakuz, T., Shin, H.M., Sima, A.K., Kitzmantl, D., Badawi, R.D., Nardo, L., Cherry, S.R., Spencer, B.A., Hacker, M., Beyer, T., 2022. Fully automated, semantic segmentation of whole-body <sup>18</sup>F-FDG PET/CT images based on data-centric artificial intelligence. *J. Nucl. Med.* 63 (12), 1941. <http://dx.doi.org/10.2967/jnumed.122.264063>, URL <http://jnm.snmjournals.org/content/63/12/1941.abstract>.
- Surveillance Research Program, National Cancer Institute, 2025. SEER\*explorer application. URL <https://seer.cancer.gov/statistics-network/explorer/>.
- Terpos, E., Mouloupos, L.A., Dimopoulos, M.A., 2011. Advances in imaging and the management of myeloma bone disease. *J. Clin. Oncol.: Off. J. Am. Soc. Clin. Oncol.* 29 (14), 1907–1915. <http://dx.doi.org/10.1200/JCO.2010.32.5449>.
- Wiegreb, S., Kopper, P., Sonabend, R., Bischl, B., Bender, A., 2024. Deep learning for survival analysis: a review. *Artif. Intell. Rev.* 57 (3), 65. <http://dx.doi.org/10.1007/s10462-023-10681-3>.
- Zamagni, E., Nanni, C., Gay, F., Pezzi, A., Patriarca, F., Bellò, M., Rambaldi, I., Tacchetti, P., Hillengass, J., Gamberi, B., Pantani, L., Magarotto, V., Versari, A., Offidani, M., Zannetti, B., Carobolante, F., Balma, M., Musto, P., Rensi, M., Mancuso, K., Dimitrakopoulou-Strauss, A., Chauviè, S., Rocchi, S., Fard, N., Marzocchi, G., Storto, G., Ghedini, P., Palumbo, A., Fanti, S., Cavo, M., 2016. 18F-FDG PET/CT focal, but not osteolytic, lesions predict the progression of smoldering myeloma to active disease. *Leukemia* 30 (2), 417–422. <http://dx.doi.org/10.1038/leu.2015.291>, URL <https://www.nature.com/articles/leu2015291>.
- Zhong, H., Huang, D., Wu, J., Chen, X., Chen, Y., Huang, C., 2023. 18F-FDG PET/CT based radiomics features improve prediction of prognosis: multiple machine learning algorithms and multimodality applications for multiple myeloma. *BMC Med. Imaging* 23 (1), 87. <http://dx.doi.org/10.1186/s12880-023-01033-2>.
- Zhu, W., Chen, Y., Nie, S., Yang, H., 2023. SAMMS: Multi-modality deep learning with the foundation model for the prediction of cancer patient survival. In: 2023 IEEE International Conference on Bioinformatics and Biomedicine. BIBM, pp. 3662–3668. <http://dx.doi.org/10.1109/BIBM58861.2023.10385661>, URL <https://ieeexplore.ieee.org/document/10385661>.