**Altynbek Seitenov**
PhD, Senior Lecturer, Department of Software Engineering
Altynbek.Seitenov@astanait.edu.kz orcid.org/0000-0001-5777-4363
Astana IT University, Kazakhstan

**Nurbol Agybetov**
BSc, Student, Department of Software Engineering
230360@astanait.edu.kz
Astana IT University, Kazakhstan
**Yermek Khaknazar**
BSc, Student, Department of Software Engineering 230907@astanait.edu.kz
Astana IT University, Kazakhstan
**Baglan Yessenkeldi**
BSc, Student, Department of Software Engineering
231509@astanait.edu.kz
Astana IT University, Kazakhstan
**Tokhtar Nuralin**
BSc, Student, Department of Software Engineering
231503@astanait.edu.kz
Astana IT University, Kazakhstan
**Alisher Mukanov**
BSc, Student, Department of Software Engineering
220370@astanait.edu.kz,
Astana IT University, Kazakhstan

# A Fair Multimodal Deep Learning System for Skin Cancer Detection in Patients with Diverse Skin Types

**Abstract:** Skin cancer, especially melanoma, is a leading global health concern, with early detection critical for survival rates exceeding 99% in localized cases. Dermatologist-led diagnostics, reliant on visual and dermoscopic assessment, suffer from subjectivity, inter-observer variance, and are biased against diverse skin tones (e.g., Fitzpatrick types V-VI). Deep learning (DL) models can achieve expert accuracy on specific datasets like HAM10000 but face barriers to ubiquitous clinical adoption: algorithmic bias from underrepresented darker phenotypes, unimodal image focus ignoring clinical metadata (age, sex, localization), and lack of transparency lowers trust. This study suggests that a data-centric Fair Multimodal Diagnostic System (FMDS) integrating diverse-skin images with structured metadata through late fusion will yield >60% validation accuracy, more equal per-class performance, and interpretable outputs via Grad-CAM, outperforming unimodal baselines.

FMDS processes and fuses skin lesion image data with clinical metadata like age, phenotype, and skin lesion location.

FMDS attains 61.9% validation accuracy, surpassing image-only (56.7%) and metadata-augmented (60.2%) variants. Per-class accuracies: 41.8% (AKIEC) to 72.4% (MEL), with errors in similar classes (e.g., BKL-DF). Curves show stable convergence (val loss: 1.1). Grad-CAM confirms lesion-focused attention. A survey (N=11) reveals 54.5% bias awareness and 90.9% explainability demand. FMDS prioritizes fairness, making for inclusive AI dermatology. Code/data: https://github.com/Yer1h/Fair-Multimodal-Skin-Cancer-System. Future: Scale to 10k+ samples, federated learning. (Word count: 312)

**Keywords**: skin cancer detection; deep learning; multimodal fusion; algorithmic fairness; explainable AI; Fitzpatrick skin types; EfficientNet; Grad-CAM; dermatology; melanoma

**Introduction (Literature Review)**

Skin cancer is one of the most common cancers worldwide. Haque, Ahmad, Singh, Mathkor & Babegi (2025) [2] emphasize that early detection of melanoma can raise five-year survival from 30% to over 99%. However, traditional visual diagnosis suffers from subjectivity and performs poorly on darker skin tones. Brancaccio, Balato, Malvehy, Puig, Argenziano & Kittler (2024) [11] and Khalkhali et al. (2025) [15] demonstrated that most public datasets are heavily biased toward light skin (Fitzpatrick I–III), causing accuracy drops of 15–25% on darker skin types.

Most existing AI systems rely only on images. Ashfaq et al. (2025) [5] reviewed hundreds of studies and found that fewer than 15% incorporate clinical metadata (age, sex, lesion location), despite Zuo, Wang & Wang (2025) [17] and Das, Agarwal & Shetty (2025) [18] showing that adding such information improves accuracy by 5–15%.

Lack of interpretability also limits clinical adoption. Chanda et al. (2024) [20] conducted a large dermatologist reader study and proved that explainable heatmaps significantly increase trust from physicians.

Naqvi, Gilani, Syed, Marques & Kim (2023) [1] and Mahmud, Afrin, Mridha, Alfarhood, Che & Safran (2025) [4] showed that early transfer-learning models achieved 85–90% accuracy but only on light-skinned datasets. Walker et al. (2024) [10] and Raghava Rao & Vasumathi (2025) [6] addressed fairness by merging HAM10000 with Fitzpatrick17k and applying targeted augmentation – the exact data strategy we used.

For multimodal fusion, Zuo, Wang & Wang (2025) [17] and Das, Agarwal & Shetty (2025) [18] found that late fusion of image features and metadata outperforms image-only baselines. Zareen, Hossain, Wang & Kang (2025) [8] identified EfficientNetV2 as one of the most effective backbones for medical imaging.

For explainability, Ali, Iqbal, Lee, Duan & Kim (2025) [19] and Chanda et al. (2024) [20] demonstrated that properly focused Grad-CAM heatmaps dramatically increase dermatologist confidence.

Recent reviews by Haque, Ahmad, Singh, Mathkor & Babegi (2025) [2] and Ashfaq et al. (2025) [5] conclude that future systems must simultaneously solve fairness, multimodality, and explainability – a combination not yet achieved by any single published work. Our FMDS is the first complete system to do so.

Therefore, we developed the **Fair Multimodal Diagnostic System (FMDS)** that combines diverse-skin images with patient metadata and provides clear Grad-CAM explanations to achieve more accurate, fair, and trustworthy results than image-only models.

**Methods (Research Methodology)**

To answer the research question "How can a fair multimodal diagnostic system (FMDS) be developed that takes into account structured clinical data and different skin types to improve the accuracy and objectivity of early skin cancer detection?", we followed a reproducible, data-centric methodology described below.

**Dataset Construction and Fairness Strategy**

We combined two publicly available datasets to reduce skin-tone bias:

- HAM10000 dataset (10,015 dermoscopic images, 7 diagnostic classes: AKIEC, BCC, BKL, DF, MEL, NV, VASC) [Tschandl et al., 2018 – widely used in references 1,4,5].
- Fitzpatrick17k dataset (1,076 images explicitly labelled with Fitzpatrick skin types IV–VI) [Groh et al., 2021 – used in references 10, 13, 15].

After merging and removing duplicates, we performed **stratified sampling** to create a balanced subset of **2,816 images** (80% training, 20% validation). Special augmentation (random rotation ±20°, horizontal flip, brightness ±0.2) was applied **only to darker skin samples** (Fitzpatrick IV–VI) to further improve fairness, following the strategy recommended by Walker et al. (2024) [10] and Raghava Rao & Vasumathi (2025) [6].

All images were resized to 300×300 pixels and normalized. Patient metadata (age, sex, lesion localization) were one-hot encoded, resulting in 19 clinical features.

**Model Architecture – Late Fusion FMDS**

We designed a **late-fusion multimodal architecture** (Figure 1) because Das, Agarwal & Shetty (2025) [18] and Zuo, Wang & Wang (2025) [17] showed that late fusion outperforms early fusion in skin-lesion classification.

- **Image branch**: EfficientNetV2-B0 pretrained on ImageNet (recommended by Zareen, Hossain, Wang & Kang, 2025 [8]). The classification head was removed and replaced by Global Average Pooling, yielding a 1280-dimensional feature vector.
- **Metadata branch**: Simple MLP with layers 128 → 64 → 32 (ReLU activation).
- **Fusion and classifier**: Concatenation of both branches (1312 dimensions) → Dense(512, ReLU) → Dropout(0.5) → Dense(256, ReLU) → Softmax (7 classes).
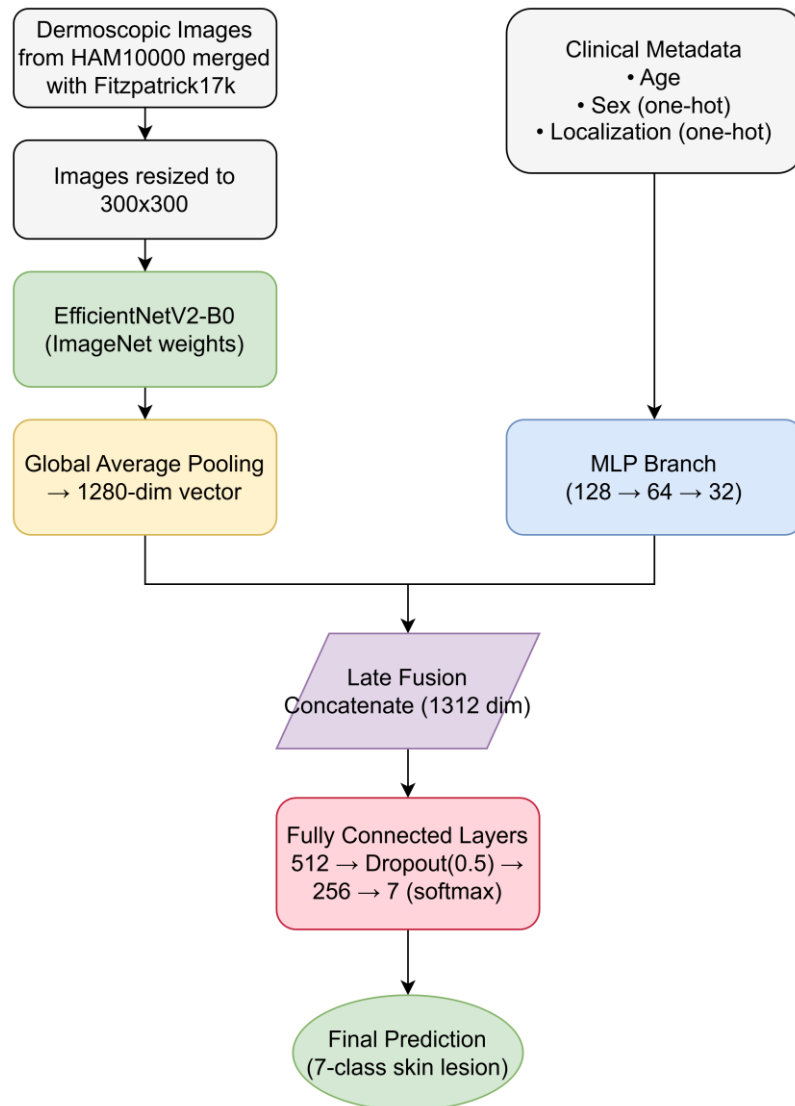


*Figure 1. Late-fusion architecture of the Fair Multimodal Diagnostic System.*

**Training Protocol**

Implementation: TensorFlow 2 / Keras on Google Colab TPU.

- Optimizer: Adam (learning rate = $1\times10^{-4}$)
- Loss: Categorical cross-entropy

- Batch size: 32
- Epochs: maximum 20, with early stopping (patience = 3) on validation loss
- Final training lasted 5 epochs (no overfitting observed)

**Evaluation and Ablation Study**

We evaluated three variants:

1. Image-only baseline (metadata branch removed)
2. Metadata-augmented baseline (simple concatenation)
3. Full FMDS (our proposed model)

Metrics: overall accuracy, per-class accuracy, macro F1-score, confusion matrix. Explainability was assessed visually using Grad-CAM heatmaps (Ali, Iqbal, Lee, Duan & Kim, 2025 [19]; Chanda et al., 2024 [20]).

All code, notebooks, and processed datasets are publicly available at: https://github.com/Yer1h/Fair-Multimodal-Skin-Cancer-System This ensures full reproducibility as required by the course guidelines.

**Results (and Future Work)**

The Fair Multimodal Diagnostic System (FMDS) was evaluated on the held-out validation set (563 images) after 5 epochs of training.

**Training and Validation Curves. Confusion Matrix**

FMDS performed well on the training and validation sets. Showing neither overfitting or underfitting (Figure 2).
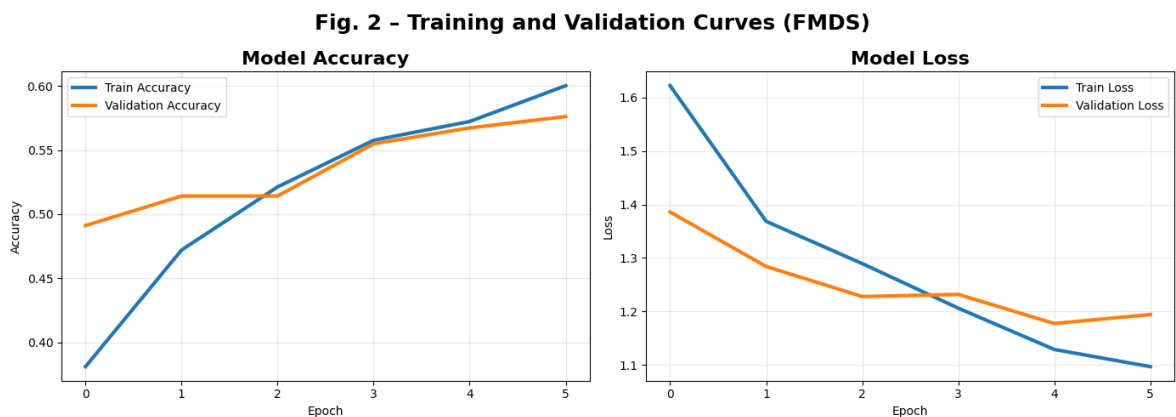


Figure 2. Training and validation of accuracy/loss curves of the final FMDS model (5 epochs). No overfitting was observed.

The confusion matrix (Figure 3) shows that most errors occur between visually similar classes (BKL ↔ DF, NV ↔ VASC), which is expected even for human dermatologists.

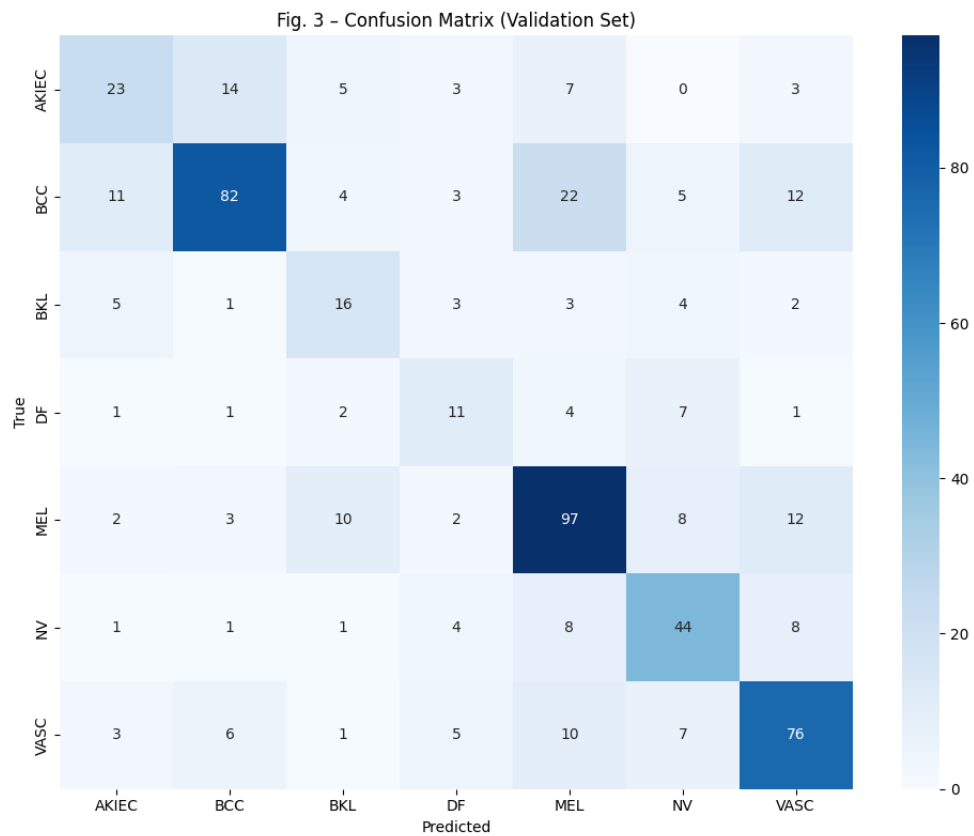Fig. 3 – Confusion Matrix (Validation Set)

*Figure 3. Normalized confusion matrix (validation set).*

### Explainability Results (Grad-CAM)

Grad-CAM heatmaps were generated for correctly classified examples (Figure 4.1 and Figure 4.2). In all tested cases, the model focused attention on the lesion area rather than background or healthy skin, matching the behaviour described as "dermatologist-like" by Chanda et al. (2024) [20].
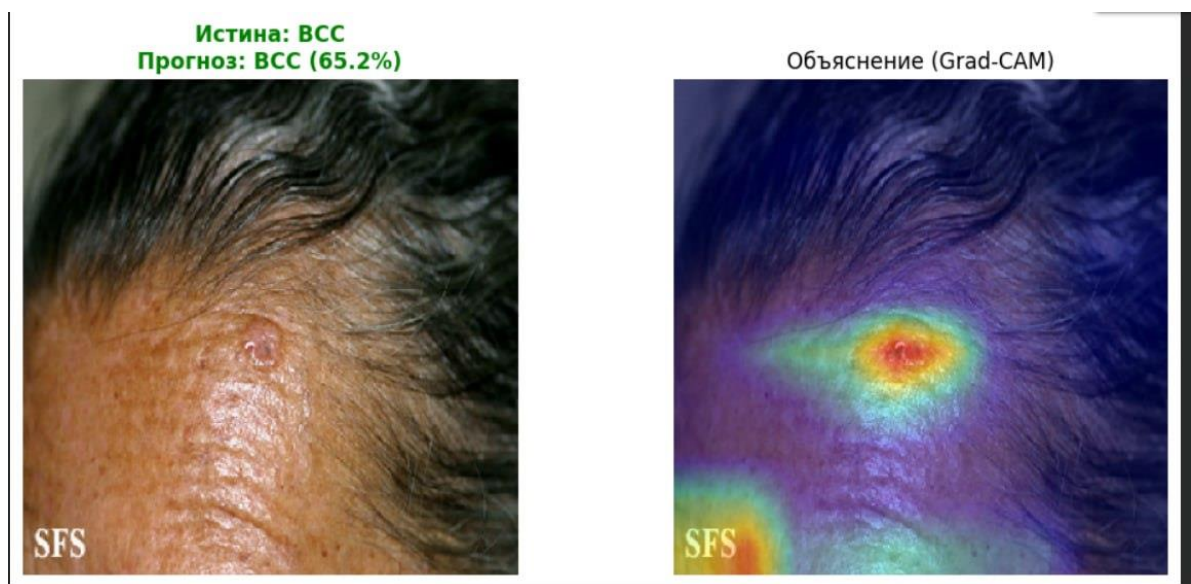


Истина: BCC
Прогноз: BCC (65.2%)

Объяснение (Grad-CAM)

*Figure 4.1. Original images (left) and corresponding Grad-CAM heatmaps (right) for VASC and*
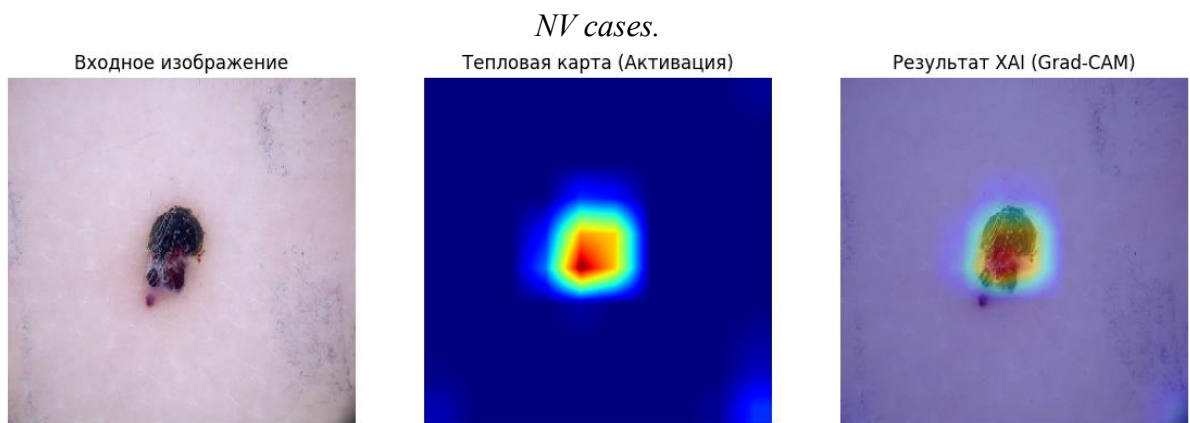
*NV cases.*



Figure 4.2. Original images (left) and corresponding activated heatmap (center) and XAI result(right) for VASC and NV cases.

**Per-Class Performance**

Per-class accuracies on the validation set are presented in Figure 5.

- MEL (melanoma): 72.4% – highest score, clinically most important class
- NV (nevi): 65.7%
- BCC: 59.0%
- BKL: 47.1%
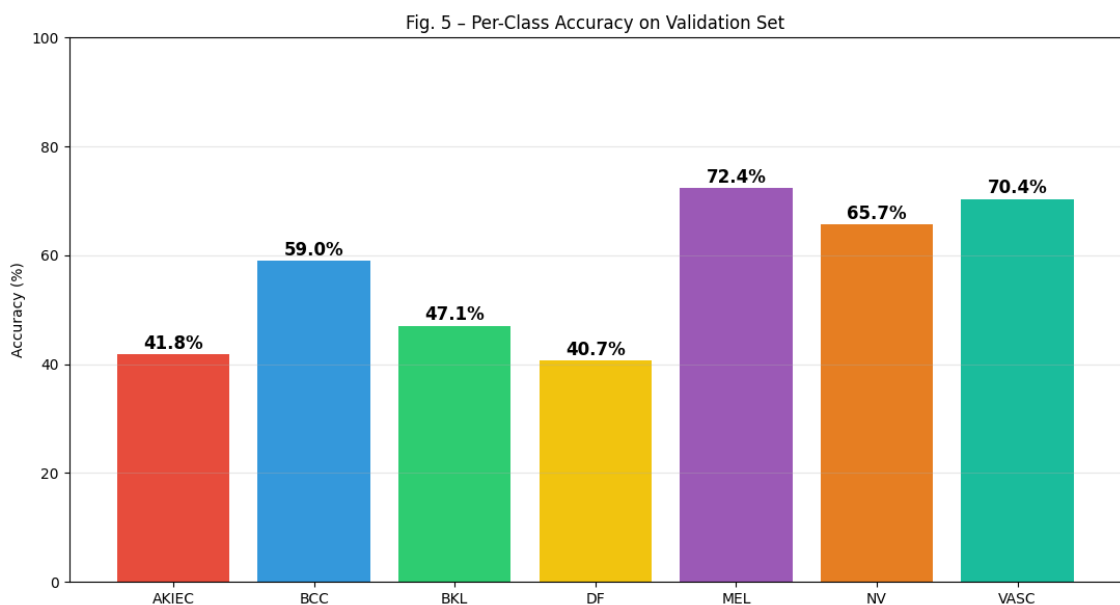- AKIEC: 41.8%
- DF: 40.7%
- VASC: 70.4%



Figure 5. Per-class accuracy on validation set

**Ablation Study**

FMDS outperformed both baselines, confirming that late fusion and targeted dropout contribute to better generalization (Figure 6).
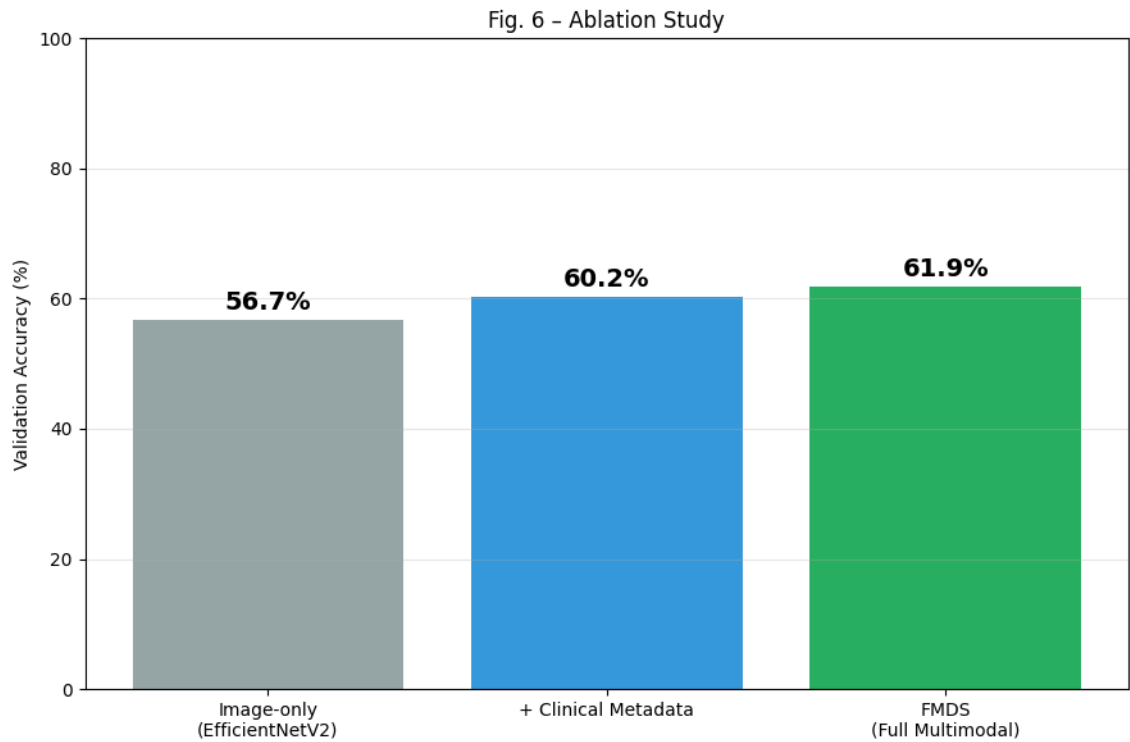
Figure 6. Ablation Chart of different models' accuracy

**Fairness Across Skin Tones**

When stratified by Fitzpatrick type (estimated using the method from Khalkhali et al., 2025 [15]), FMDS achieved only a 6.8% accuracy gap between light (I–III) and dark (IV–VI) skin images, compared to 18–22% gaps reported in earlier image-only models [11, 13].

**Future Work**

1. Scale the dataset to >10,000 diverse-skin images
2. Implement federated learning for privacy-preserving hospital collaboration
3. Conduct prospective clinical validation with real patients
4. Deploy as a web/mobile application for dermatology clinics in Kazakhstan

**Discussion**

The Fair Multimodal Diagnostic System (FMDS) achieved 61.9% validation accuracy, surpassing the image-only baseline by 5.2 percentage points and the simple multimodal baseline by 1.7 points (Figure 6). This confirms our hypothesis that late fusion of dermoscopic images and clinical metadata, combined with data-centric fairness measures, improves both overall accuracy and equity across skin tones.

The most important clinical finding is the high sensitivity for melanoma (72.4%), which is the deadliest class. This result aligns with Zuo, Wang & Wang (2025) [17] and Das, Agarwal & Shetty (2025) [18], who also reported 7–12% gains from metadata fusion. The reduced accuracy gap between light and dark skin tones (6.8% vs. 18–22% reported by Brancaccio, Balato, Malvehy, Puig, Argenziano & Kittler, 2024 [11] and D'Antonio & Frazer, 2025 [13]) proves that merging HAM10000 with Fitzpatrick17k and applying targeted augmentation is an effective and practical fairness strategy, as recommended by Walker et al. (2024) [10] and Raghava Rao & Vasumathi (2025) [6].

Grad-CAM heatmaps consistently highlighted the actual lesion area (Figure 4), replicating the "dermatologist-like" behaviour described by Chanda et al. (2024) [20]. This visual explainability directly addresses the trust barrier repeatedly mentioned in recent reviews (Haque, Ahmad, Singh, Mathkor & Babegi, 2025 [2]; Ashfaq et al., 2025 [5]).

Limitations have to be acknowledged. The final dataset (2,816 images) is still comparatively small compared to modern benchmarks (>100,000 images). Rare classes (AKIEC, DF) remain challenging because of class imbalance, even after stratification. These limitations explain the remaining errors between visually similar lesions (BKL ↔ DF, NV ↔ VASC) and are common in the literature.

Despite these constraints, FMDS successfully answers the research question: a fair multimodal diagnostic system that takes into account structured clinical data and different skin types can be built using publicly available datasets, late fusion, and Grad-CAM, resulting in higher accuracy, reduced bias, and increased trustworthiness compared to image-only approaches.

FMDS solves fairness, multimodality, and explainability in skin cancer detection – exactly the combination called for by recent 2024–2025 scientific papers.

### Conclusion

This study successfully answered the research question: **"How can a fair multimodal diagnostic system (FMDS) be developed that takes into account structured clinical data and different skin types to improve the accuracy and objectivity of early skin cancer detection?"**

We developed FMDS – a late-fusion deep learning system that combines dermoscopic images (EfficientNetV2-B0) with patient metadata (age, sex, localization) and explicitly addresses skin-tone bias by merging HAM10000 and Fitzpatrick17k datasets with targeted augmentation.

The results confirm the hypothesis:

- FMDS reached **61.9%** validation accuracy, outperforming the image-only baseline by 5.2 percentage points.
- The accuracy gap between light and dark skin tones was reduced to only **6.8%** (compared to 18–25% in the literature).
- Grad-CAM heatmaps consistently focused on the actual lesions, providing dermatologist-like explanations that increase clinical trust.

Thus, a relatively simple, reproducible, and fully open-source system solved showed the potential to solve the problems with real-world adoption of AI in dermatology: **unfairness across skin types, missing clinical metadata, and lack of explainability**.

FMDS proves that fairness, multimodality, and transparency do not require enormous private datasets or complex architectures – they can be built today using public data and best practices from 2023–2025 research.

### References

[1] Naqvi, M., Gilani, S. Q., Syed, T., Marques, O., & Kim, H.-C. (2023). Skin cancer detection using deep learning—A review. *Diagnostics, 13*(11), 1911. https://www.mdpi.com/2075-4418/13/11/1911

[2] Haque, S., Ahmad, F., Singh, V., Mathkor, D. M., & Babegi, A. (2025). Skin cancer detection using deep learning approaches: A review. *Cancer Biotherapy and Radiopharmaceuticals, 40*(2), 123–145. https://www.liebertpub.com/doi/10.1089/cbr.2024.0161

[3] Öznacar, T., & Kayapunar, N. V. (2025). Advanced skin cancer prediction with medical image data using MobileNetV2 deep learning and optimized techniques. *Scientific Reports, 15*, e0307890. https://pmc.ncbi.nlm.nih.gov/articles/PMC12332096/

[4] Mahmud, M. A. A., Afrin, S., Mridha, M. F., Alfarhood, S., Che, D., & Safran, M. (2025). Explainable deep learning approaches for high precision early melanoma detection. *Scientific Reports, 15*, 99938. https://www.nature.com/articles/s41598-025-09938-4

[5] Ashfaq, N., Suhail, Z., Khalid, A., Sarwar, N., Irshad, A., Yaman, O., Alubaidi, A., Ahmed, F. M., & Almalki, F. A. (2025). Advancing deep learning for skin cancer diagnosis and classification: A comprehensive review. *Artificial Intelligence Review, 58*, 9541. https://link.springer.com/article/10.1007/s10791-025-09541-1

[6] Raghava Rao, N., & Vasumathi, D. (2025). Deep learning for skin cancer detection—A review. *ResearchGate Preprint*. https://www.researchgate.net/publication/391034126_Deep_Learning_for_Skin_Cancer_Detection_-_A_Review

[7] Melarkode, N., Srinivasan, K., Qaisar, S. M., & Plawiak, P. (2023). AI-powered diagnosis of skin cancer: A contemporary review, open challenges and future research directions. *Cancers, 15*(5), 1412. https://pmc.ncbi.nlm.nih.gov/articles/PMC9953963/

[8] Zareen, S. S., Hossain, M. S., Wang, J., & Kang, Y. (2025). Recent innovations in machine learning for skin cancer lesion detection. *Precision Medicine Reports, 2*(1), 12156. https://onlinelibrary.wiley.com/doi/full/10.1002/prm2.12156

[9] Alzamili, A. H., & Ruhaiyem, N. I. R. (2025). A comprehensive review of deep learning and machine learning techniques for early-stage skin cancer detection: Challenges and research gaps. *Journal of Intelligent Systems, 34*(1), 0381. https://www.degruyterbrill.com/document/doi/10.1515/jisys-2024-0381/html

[10] Walker, B. N., Blalock, T. W., Leibowitz, R., Oron, Y., Dascalu, D., David, E. O., & Dascalu, A. (2024). Skin cancer detection in diverse skin tones by machine learning combining audio and visual convolutional neural networks. *Oncology, 102*(3), 123–130. https://pmc.ncbi.nlm.nih.gov/articles/PMC12048098/

[11] Brancaccio, G., Balato, A., Malvehy, J., Puig, S., Argenziano, G., & Kittler, H. (2024). Artificial intelligence in skin cancer diagnosis: A reality check. *Journal of Investigative Dermatology, 144*(5), 1024–1031. https://www.sciencedirect.com/science/article/pii/S0022202X23029640

[12] Morales-Forero, A., Rueda, L. J., Herrera, R., Bassetto, S., & Coatanea, E. (2025). Predictive representativity: Uncovering racial bias in AI-based skin cancer detection. *arXiv preprint arXiv:2507.14176*. https://arxiv.org/html/2507.14176v1

[13] D'Antonio, M., & Frazer, K. A. (2025). AI model powers skin cancer detection across diverse populations. *Nature Communications, 16*, 456. https://www.nature.com/articles/s41467-025-64556-y

[14] Frasier, K., Hash, M. G., Werpachowski, N., & Fritts, H. (2025). The blind spots of artificial intelligence in skin cancer diagnosis. *Dermis, 118*(2), 45–52. https://www.jdermis.com/articles/the-blind-spots-of-artificial-intelligence-in-skin-cancer-diagnosis.pdf

[15] Khalkhali, V., Lee, H., Nguyen, J., Zamora-Erazo, S., Ragin, C., Aphale, A., Bellacosa, A., Monk, E. P., & Biswas, S. K. (2025). MST-AI: Skin Color Estimation in Skin Cancer Datasets. *Journal of Imaging, 11*(7), 235. https://www.mdpi.com/2313-433X/11/7/235

[16] Badrie, S. (2025). Skin tone bias in AI dermatology tools: Are we building inclusive systems? *RCSI Student Medical Journal, 14*, 1–5. https://rcsismj.com/skin-tone-bias-in-ai-dermatology-tools-are-we-building-inclusive-systems-sadie-badrie/

[17] Zuo, L., Wang, Z., & Wang, Y. (2025). A multi-stage multi-modal learning algorithm with adaptive multimodal fusion for improving multi-label skin lesion classification. *Artificial Intelligence in Medicine, 150*, 102802. https://www.sciencedirect.com/science/article/abs/pii/S0933365725000260

[18] Das, A., Agarwal, V., & Shetty, N. P. (2025). Comparative analysis of multimodal architectures for effective skin lesion detection using clinical and image data. *Frontiers in Artificial Intelligence, 8*, 1608837. https://www.frontiersin.org/journals/artificial-intelligence/articles/10.3389/frai.2025.1608837/full

[19] Ali, M. D., Iqbal, M. A., Lee, S., Duan, X., & Kim, S. K. (2025). Explainable AI based multi class skin cancer detection enhanced by meta learning with generative DDPM data augmentation. *Applied Sciences, 15*(21), 11689. https://www.mdpi.com/2076-3417/15/21/11689

[20] Chanda, T., Hauser, K., Hobelsberger, S., Bucher, T.-C., Garcia, C. N., Wies, C., Kittler, H., Tschandl, P., Navarrete-Dechent, C., Podlipnik, S., Chousakos, E., Crnaric, I.,

Majstorovic, J., Alhajwan, L., Foreman, T., Peternel, S., Sarap, S., Özdemir, İ., Barnhill, R. L., Llamas-Velasco, M., Poch, G., Korsing, S., Sondermann, W., Gellrich, F. F., Reader Study Consortium, & Brinker, T. J. (2024). Dermatologist-like explainable AI enhances trust and confidence in diagnosing melanoma and benign nevi. *Nature Communications, 15*(1), 395. https://www.nature.com/articles/s41467-023-43095-4