

# 딥페이크 히트맵 분석 보고서

## 1. 분석 개요

- 모델명: MobileNetV3-Small
- 모델 유형: 한국인 전용 모델
- 분석 일시: 2025-11-05 17:34:48
- 예측 결과: 이 이미지는 Fake로 분류되었으며, 예측신뢰도는 85.68%입니다.
- 비정상적인 질감, 경계선 왜곡, 조명 불균형 등 딥페이크 생성 흔적이 감지되었습니다.
- 딥페이크 확률: 0.00%

## 2. Grad-CAM 시각화



## 3. LangChain 기반 AI 해석

Grad-CAM (Gradient-weighted Class Activation Mapping)은 딥러닝 모델의 예측 결과에 대한 시각적 해석을 제공하는 유용한 기술입니다. 주어진 정보와 Grad-CAM 히트맵을 기반으로 딥페이크 탐지 결과를 해석해보겠습니다.

### ### 모델 판단 근거 해석

1. \*\*붉은색 영역\*\*: 히트맵에서 붉은색으로 강조된 영역은 모델이 이 이미지가 딥페이크로 분류된 주된 근거를 지원하는 부분입니다. 이 색상은 모델이 해당 부분에서 중요한 특징을 발견했음을 의미합니다.

## 2. \*\*합성 혼적\*\*:

- \*\*비정상적인 질감\*\*: 붉은색 영역에서 피부 질감이 고르지 않거나 불균형하게 나타나면, 이는 인간의 자연스러운 피부 결과 차별화됩니다. 딥페이크는 종종 피부 질감을 부자연스럽게 처리하곤 하며, 이로 인해 고르지 않은 색상이나 불일치가 생깁니다.

- \*\*경계선 왜곡\*\*: 이미지의 경계선이 부자연스럽게 흐릿하거나 왜곡되어 나타나는 경우, 이는 합성이 이루어졌음을 나타냅니다. 예를 들어, 얼굴과 배경의 경계가 명확하지 않거나 비정상적으로 혼합된 패턴이 관찰될 수 있습니다.

## 3. \*\*조명 불균형\*\*:

- 모델이 특정 부위의 조명 상태가 불균형하게 나타나는 부분을 강조할 경우, 이는 종종 합성된 영역과 실제 이미지 간의 조명 차이에서 발생합니다. 자연스러운 조명은 인물의 얼굴을 부드럽게 감싸지만, 딥페이크는 특정 영역에서 조명 조건이 다르게 나타날 수 있습니다.

## ### 신뢰도 및 한계점

### 1. \*\*신뢰도\*\*:

- 예측신뢰도 85.68%는 상당히 높은 수치로, 모델이 이 이미지를 딥페이크로 잘 분류한 것을 나타냅니다. 이는 많은 참고 사례와 충분한 학습을 바탕으로 하였을 가능성이 큽니다.

- 그러나 신뢰도가 100%가 아닌 이상, 언제나 잘못된 분류 가능성이 존재합니다. 예를 들어, 자연스러운 이미지에서 조명이나 텍스처가 비정상적으로 보일 수 있는 상황이 발생할 수 있습니다.

### 2. \*\*한계점\*\*:

- Grad-CAM은 모델이 어떤 부분을 중요하게 생각하는지 시각적으로 보여줄 수 있지만, 그 해석이 항상 정확한 것은 아닙니다. 예를 들어, 실제 이미지에서 사람의 표정이나 조명 조건이 변했을 경우, 모델이 높은 신뢰도로 잘못된 결론을 내릴 수 있습니다.

- 또한, 특정한 상황에서 통계적 우연에 의해 잘못된 예측이 발생할 수 있으며, 검증되지 않은 경우 신뢰성에 대한 의문을 낳을 수 있습니다.

## ### 심층 결과

- \*\*추가적인 시각적 근거\*\*: 이미지의 일부가 포함된 인물의 표정 변화나 감정 표현의 비정상성이 나타난다면, 이는 추가적인 심층 결과로 볼 수 있습니다. 딥페이크 생성 기술은 다양한 얼굴 표정의 표현을 결합하는 데 있어 한계가 있을 수 있습니다.
- \*\*컨텍스트 이해 부족\*\*: 모델이 이미지의 전체 맥락(예: 배경, 인물 간의 상호작용 등)을 인식하지 못하는 경우, 특정 세부사항에 과도하게 의존할 수 있습니다. 이는 잘못된 긍정적 결과를 초래할 수 있습니다.

## ### 결론

모델은 딥페이크의 여러 비정상적인 특징을 기반으로 높은 신뢰도로 이미지를 분류했지만, Grad-CAM 해석을 통해 시각적으로 제시된 증거를 평가하면서, 신뢰도와 한계점을 고려하는 것이 중요한 과정입니다. 이러한 해석 및 분석은 딥페이크 탐지 기술을 더욱 발전시키고, 인식력을 높이는 데 기여할 수 있습니다.

#### 4. 결론 및 권장 조치

본 분석은 Grad-CAM 시각 주목도를 중심으로 진행되었습니다.

AI의 결과는 참고용으로 사용해야 하며, 법적 판단이나 공식 증거로 사용되지 않습니다.

결과의 신뢰도를 높이기 위해 다양한 이미지 소스로 교차 검증을 권장합니다.