

# 딥페이크 히트맵 분석 보고서

## 1. 분석 개요

- 모델명: MobileNetV3-Small
- 모델 유형: 외국인 전용 모델
- 분석 일시: 2025-11-07 10:36:54
- 예측 결과: Fake (52.49%)
- 딥페이크 확률: 18.80%

## 2. Grad-CAM 시각화



## 3. LangChain 기반 AI 해석

Grad-CAM(Gradient-weighted Class Activation Mapping)은 딥러닝 모델이 특정 클래스를 예측할 때 어떤 영역에 주목했는지를 시각적으로 표현해주는 기법입니다. 주어진 정보를 기반으로 Grad-CAM 히트맵을 해석해보겠습니다.

### ### 1. Grad-CAM 히트맵 해석

- \*\*붉은색 영역\*\*: 모델이 딥페이크를 판단하는 데 있어서 중요한 시각적 특징을 나타내는 곳입니다. 이 지역은 다음과 같은 요소들을 포함할 가능성이 높습니다:

- \*\*합성 흔적\*\*: 이미지의 경계나 곡선이 자연스러운 피부 결이나 특징과 맞지 않을 수 있습니다. 만약 붉은색 영역이 얼굴의 경계, 턱, 또는 이마 부분에 있다면, 합성된 부분에서의 불연속성이 의심스러울 수

있습니다.

- **피부 질감**: 딥페이크 모델들은 종종 비현실적인 피부 질감을 만들기 때문에, 붉은 영역에서 피부의 세부적인 질감이 다른 영역과 다르다면 그것이 모델의 판단 근거가 될 수 있습니다. 예를 들어, 과도하게 부드럽거나 과도한 노이즈가 발생할 수 있습니다.

- **조명 왜곡**: 진짜 이미지와 비교했을 때 조명이 비대칭하거나 비정상적으로 분포되어 있는 경우, 이러한 요소도 모델의 판단에 큰 영향을 미칠 수 있습니다. 예를 들어, 특정 부분에서 그림자가 부자연스럽게 보이거나 빛의 반사가 자연스럽지 않다면, 그 부분이 붉은색으로 표시될 수 있습니다.

### ### 2. 신뢰도와 한계점

- **신뢰도**: Grad-CAM은 모델의 신뢰도를 강화해주는 도구이지만, 여전히 한계가 존재합니다. 모델이 붉은색 영역에 주목했다고 해서 해당 부분이 반드시 딥페이크의 확실한 증거가 되지는 않습니다. 모델이 경향성을 가진 학습 데이터를 기반으로 판단했기 때문에, 일부 딥페이크 이미지의 경우 실제 인물의 얼굴로 판단할 수 있는 오류가 발생할 가능성이 있습니다.

- **한계점**:

- **가짜 이미지의 질**: 딥페이크의 품질에 따라서 Grad-CAM의 해석이 달라질 수 있습니다. 높은 퀄리티의 딥페이크는 더 많은 자연스러운 특징을 가질 수 있어 모델의 판단이 모호해질 수 있습니다.

- **일반화 문제**: 특정 데이터셋에서 학습한 모델은 새로운 이미지에 대해 잘 일반화하지 못할 수 있습니다. 결과적으로, 새로운 유형의 딥페이크가 나타날 경우, 모델이 해당 이미지를 올바르게 분류하지 못하는 경우도 있습니다.

### ### 3. 추가적인 심층 결과

- **특징 분리에 대한 이해**: 모델이 특정 특징을 추출하는 과정을 이해함으로써 내부 동작 방식을 더욱 깊이 있게 분석할 수 있습니다. 예를 들어, 모델이 어떤 피부 톤이나 주름, 또는 표정을 인식하는 방식은 특정 딥페이크 기술에 따라 변화할 수 있습니다.

- **다양한 각도와 조명**: 모델은 다양한 각도와 조명 조건에서의 인물의 신뢰도를 평가하여, 예를 들어 비대칭적인 조명이나 각도로 인한 불일치를 감지할 수 있습니다. 이는 딥페이크 탐지의 중요한 요소가 됩니다.

- **변별력 있는 피처**: 고주파, 저주파 정보의 차이와 같이, 이미지의 세밀한 변별력을 분석하는 것이 가능할 경우, 더 깊이 있는 판단이 가능합니다. 이 경우에는 심층적으로 각 층의 출력을 분석하여 특징적인 패턴을 찾아낼 수 있습니다.

이와 같이 Grad-CAM 히트맵을 통해 모델의 판단 근거를 파악하는 것은 모델을 이해하고 성능을 개선하는데 유용하지만, 깊은 이해와 추가적인 분석 없이 결과를 단순하게 해석하는 것은 위험할 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.