

| 사용자 맞춤형 금지어 필터링을 적용한 실시간 온라인 채팅 시스템

숭실대학교 소프트웨어학부

팀 : 창호공들

목 차

1
개발 배경 & 개발 목적

2
관련 연구

3
사용자 맞춤형 금지어 필터링 모델

4
실험

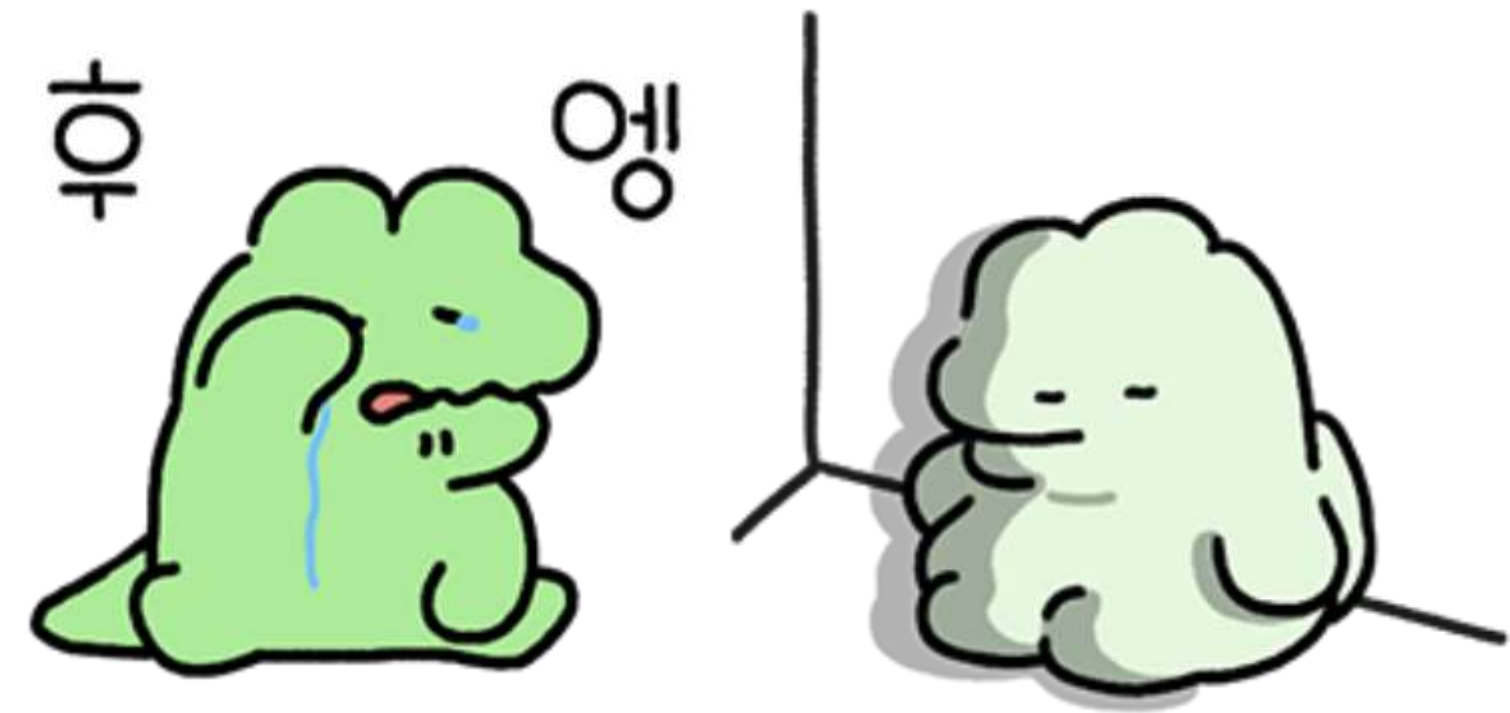
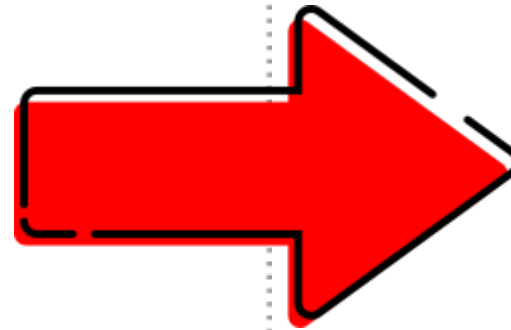
5
실험 결과 분석

6
결론 및 향후 연구 방향

개발 배경

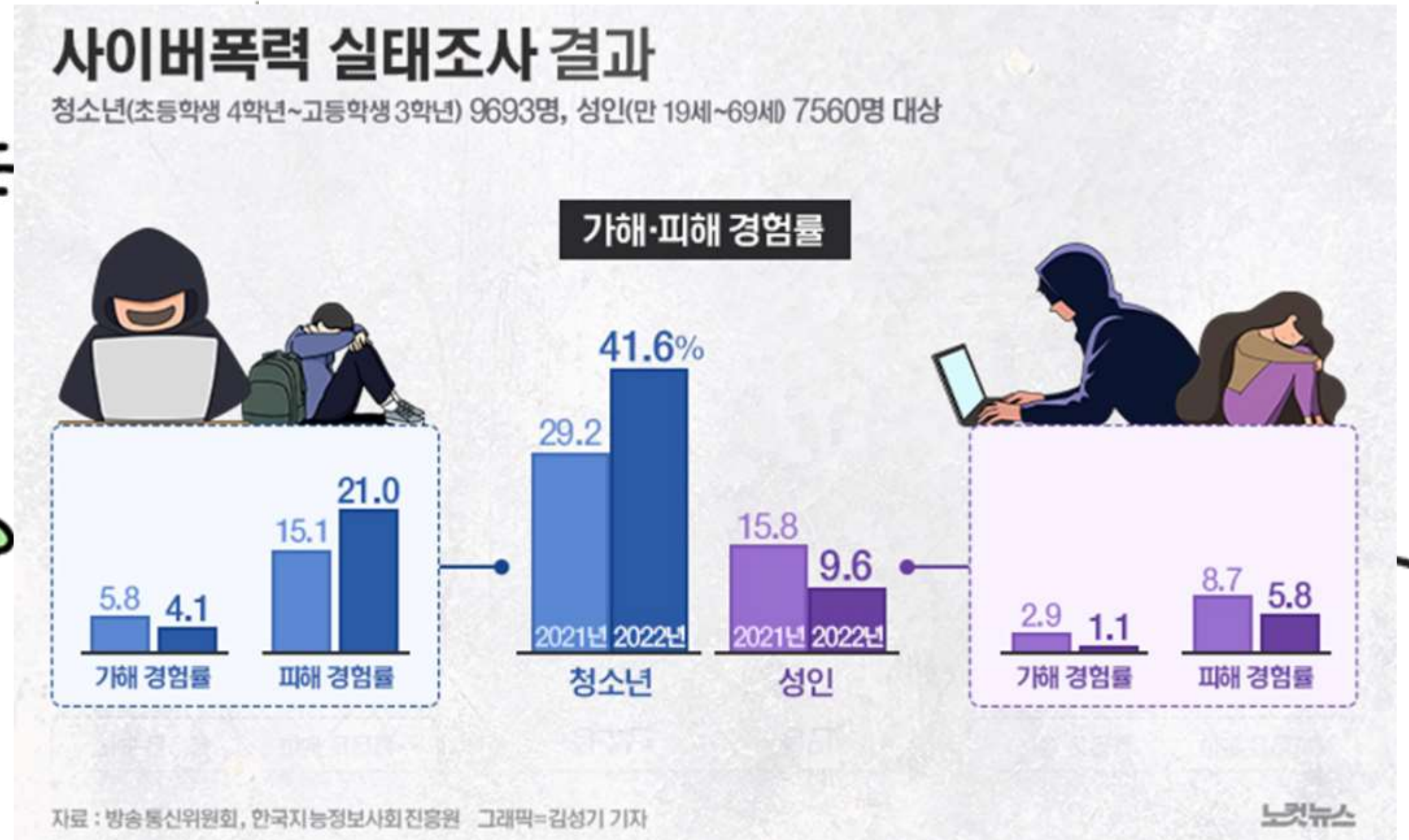
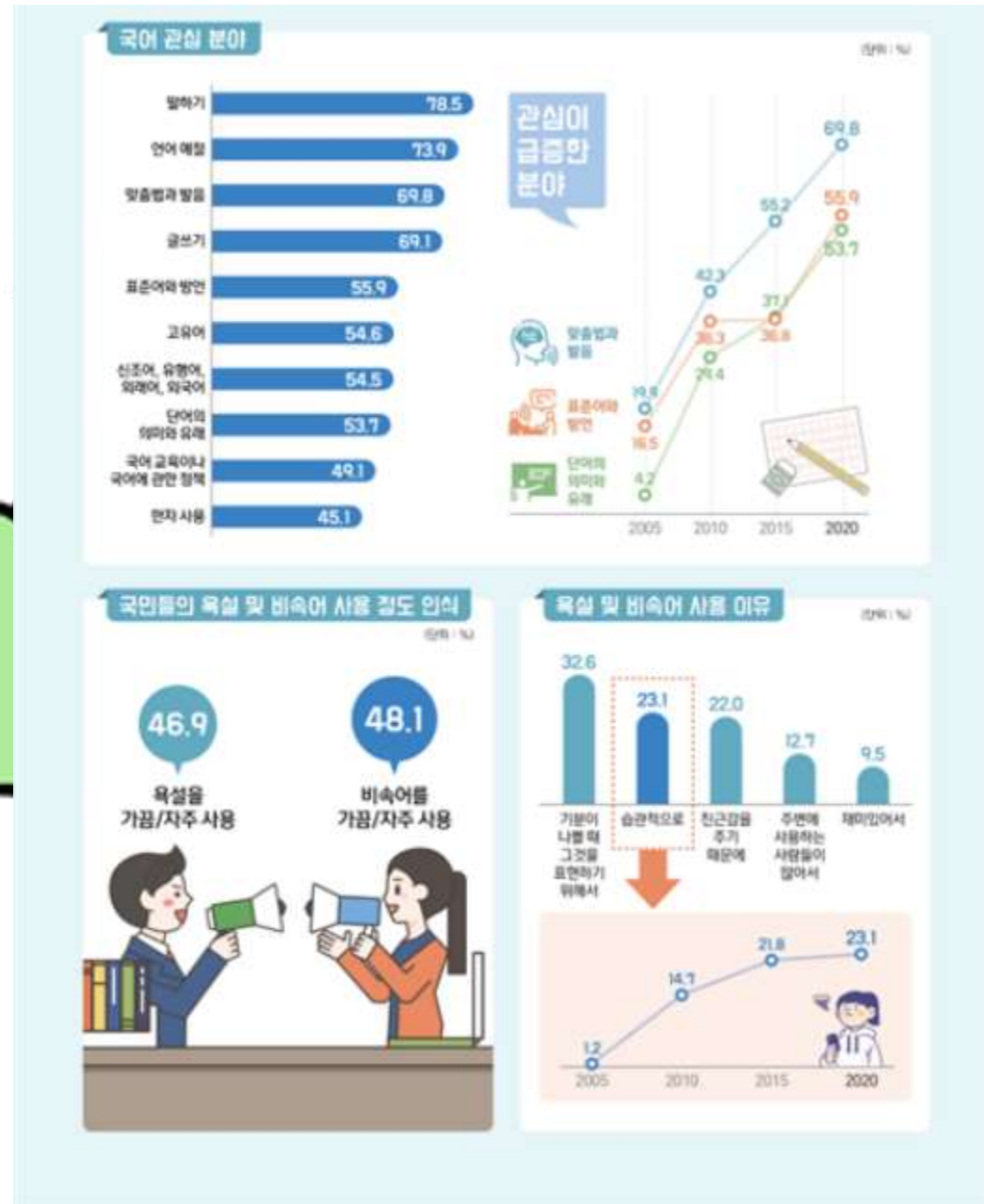


욕설 ↑



사회적 문제 ↑

개발 배경



개발 목표



- 욕설, 비속어 금지어 목록 생성
 - 금지어 목록의 욕설, 비속어 실시간 필터링
- > 온라인 커뮤니케이션 환경 개선 & 언어 사용에 대한 인식 개선

관련 연구

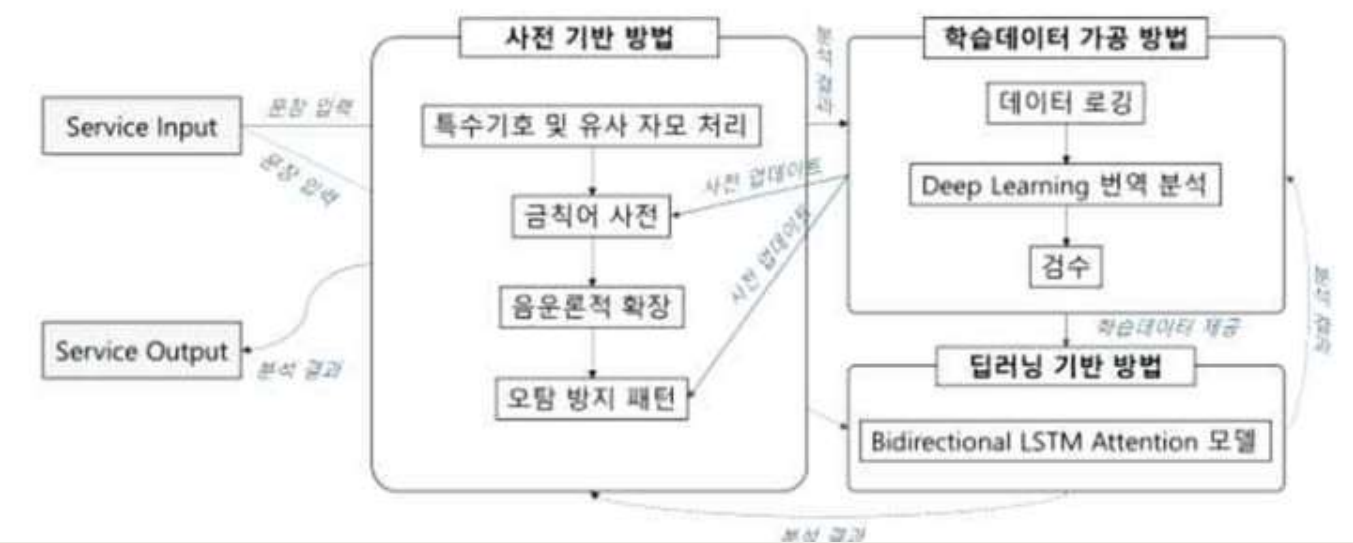
1. 정규 표현식 기반 필터링
2. 사전 기반 키워드 필터링
3. 딥러닝 기반 메시지 필터링

1. 정규 표현식 기반 필터링

- 미리 정의된 패턴을 기반으로 특정 단어 탐지
- 장점 : 한 개의 단어를 가지고 유사한 패턴의 단어를 필터링 가능
- 단점 : "눔" -> "ㄴ ㄴ ㄴ"과 같이 변형 시 필터링 불가

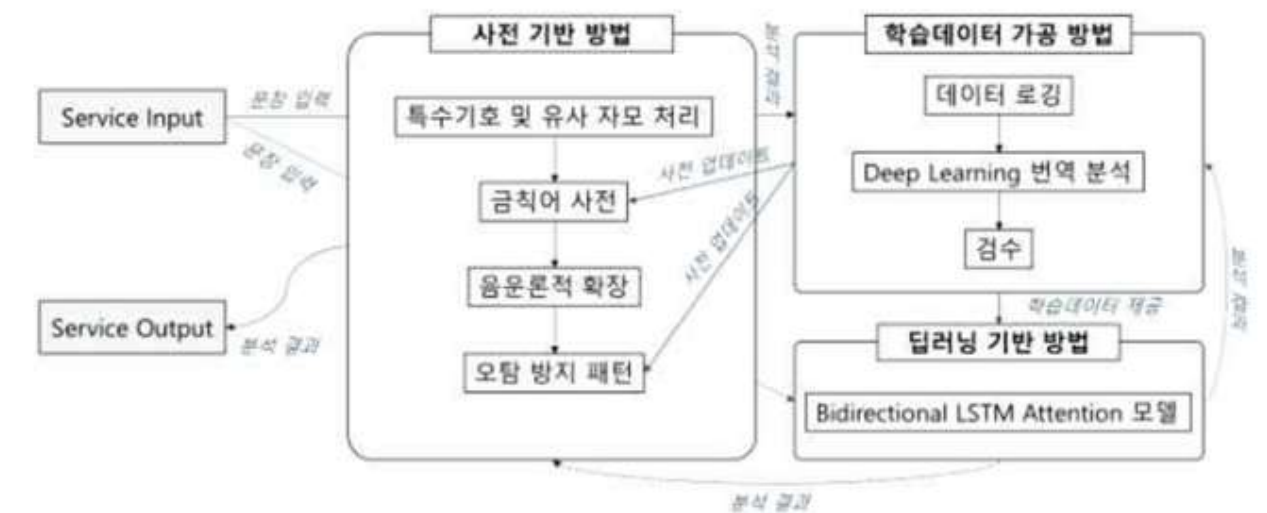
2. 사전 기반 키워드 필터링

- 금지어 목록을 설정, 메시지에서 해당 단어가 포함되는지 검사
- 장점 : 저장된 금지어에 대한 관리가 용이함
- 단점 : 추가되지 않은 단어에 대한 필터링이 불가함

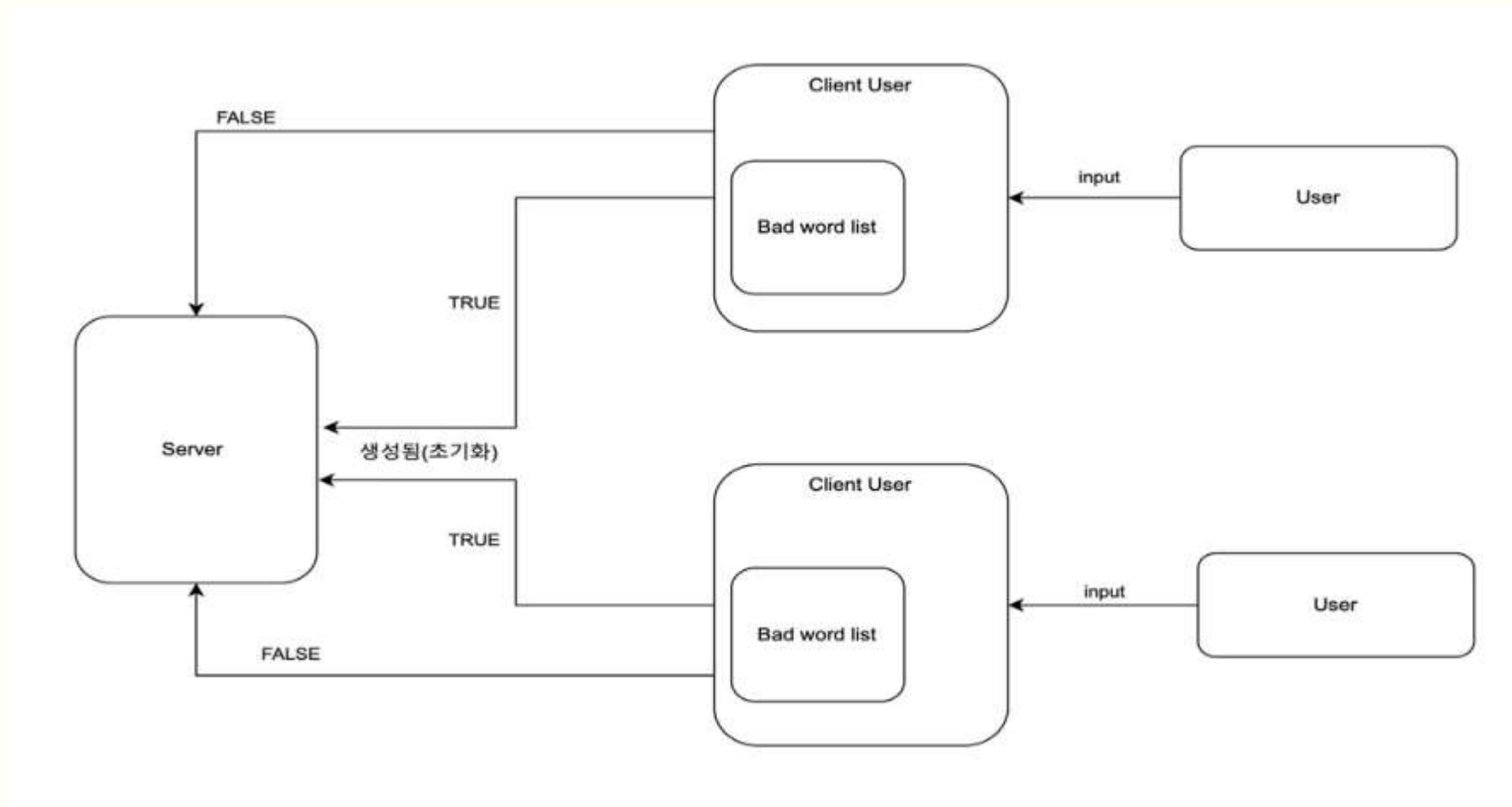


3. 딥러닝 기반 메시지 필터링

- 학습집합을 반자동적으로 확보, 딥러닝 모델이 지속적으로 학습 가능한 메시지 필터링 구조를 제안하여 해당 메시지 필터를 통해 필터링
- 장점 : 앞선 방식의 단점을 보완 가능



사용자 맞춤형 금지어 필터링 모델



사용한 필터링 모델 : 사전 기반 키워드 필터링

필터링 과정 :

1. 개인화된 금지어 목록 관리

2. 메시지 입력 및 수신

3. 단어 탐지 및 대체

4. 메시지 출력 및 전송

실험 절차

채팅 서버 작동 및 클라이언트 서버 연결

채팅 서버를 실행하고
클라이언트를 서버에
접촉하여 실시간 채팅
환경 구축



채팅 테스트

실시간으로 채팅이
잘 보내지는지 확인



금지어 데이터셋 추가

필터링을 하고자하는
언어를 금지어 목록에 추가

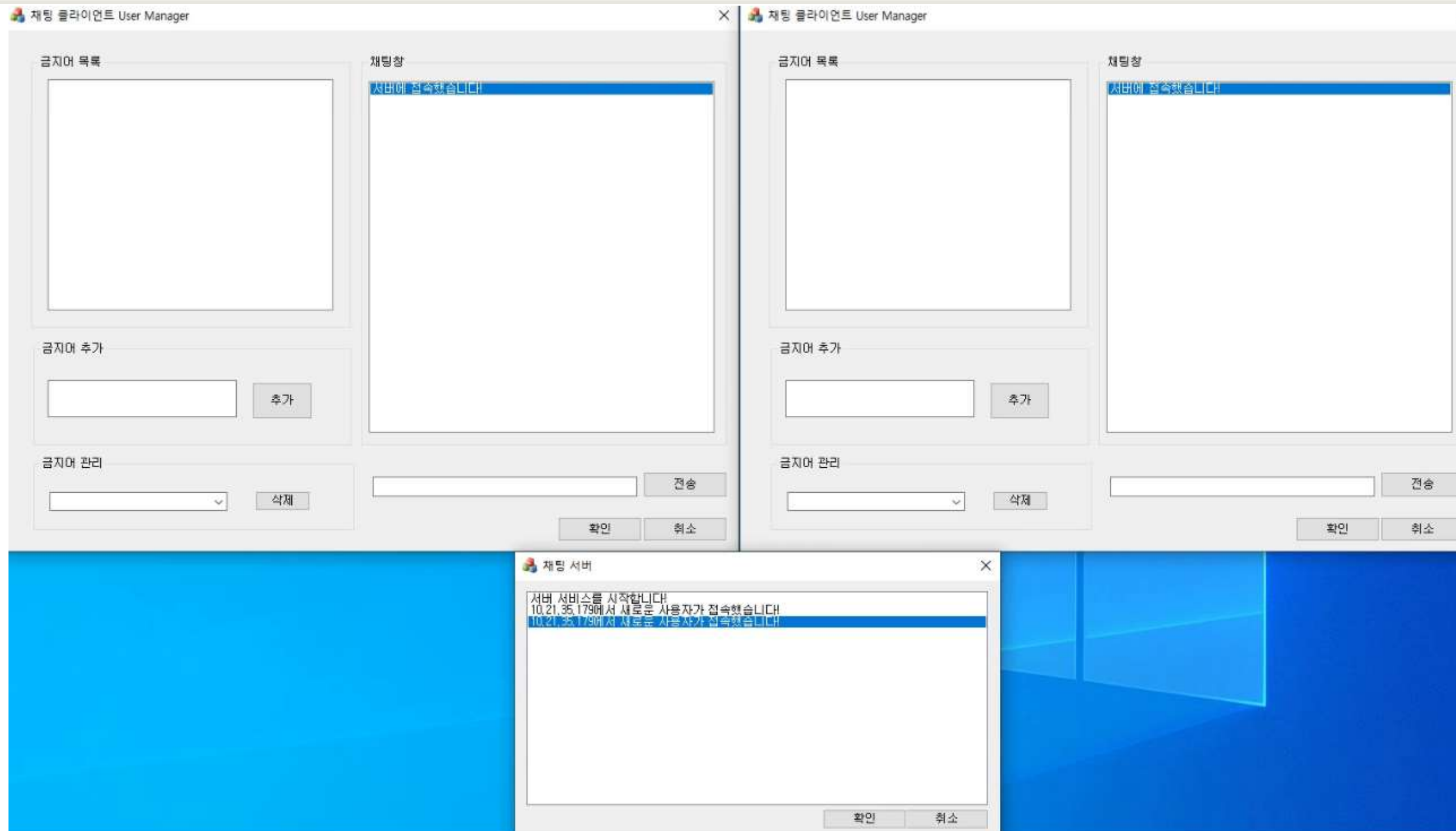


필터링 여부 확인

금지어 목록에 추가된 문자가
제대로 필터링되는지 확인



실험 및 실험 결과

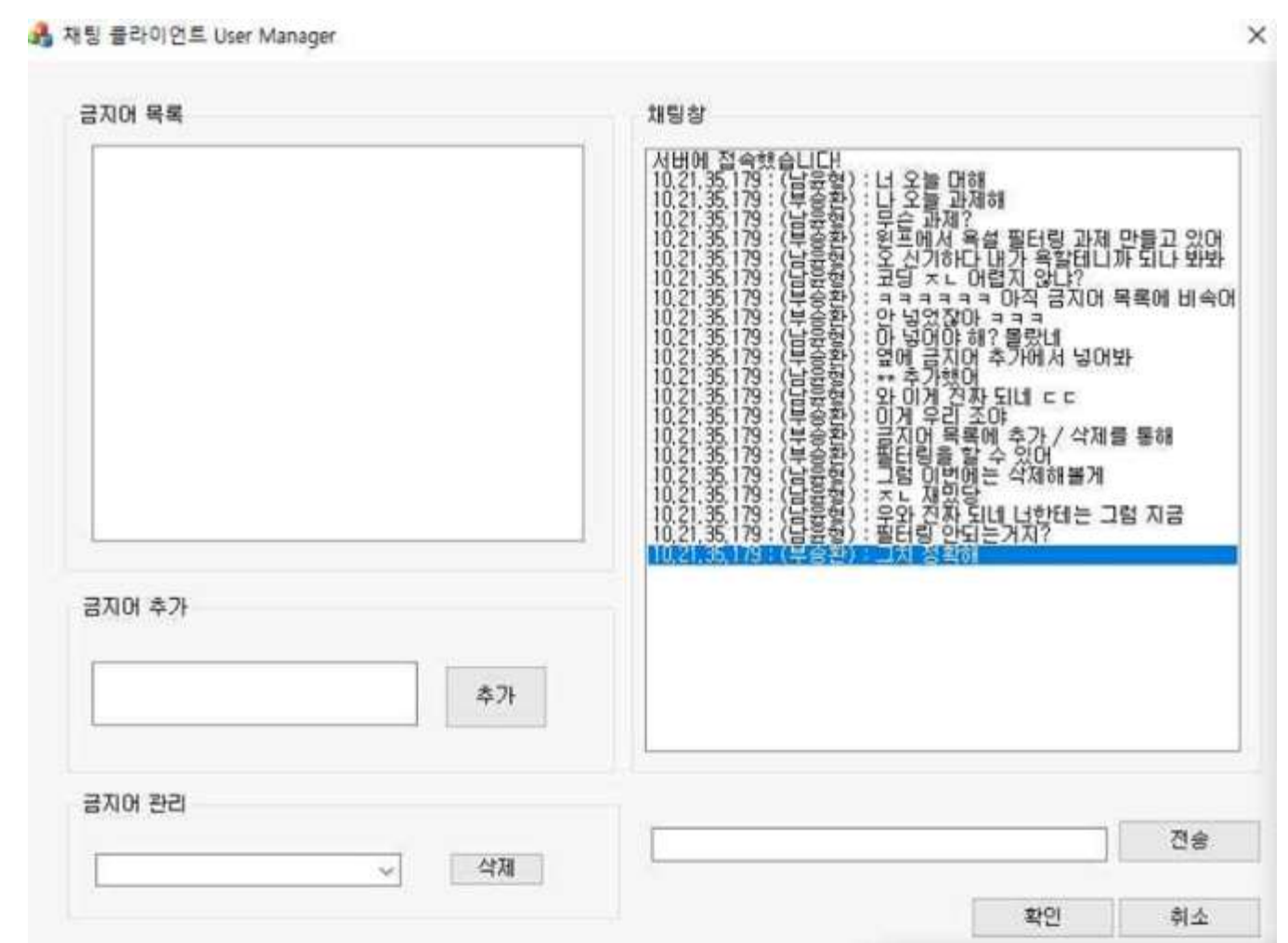
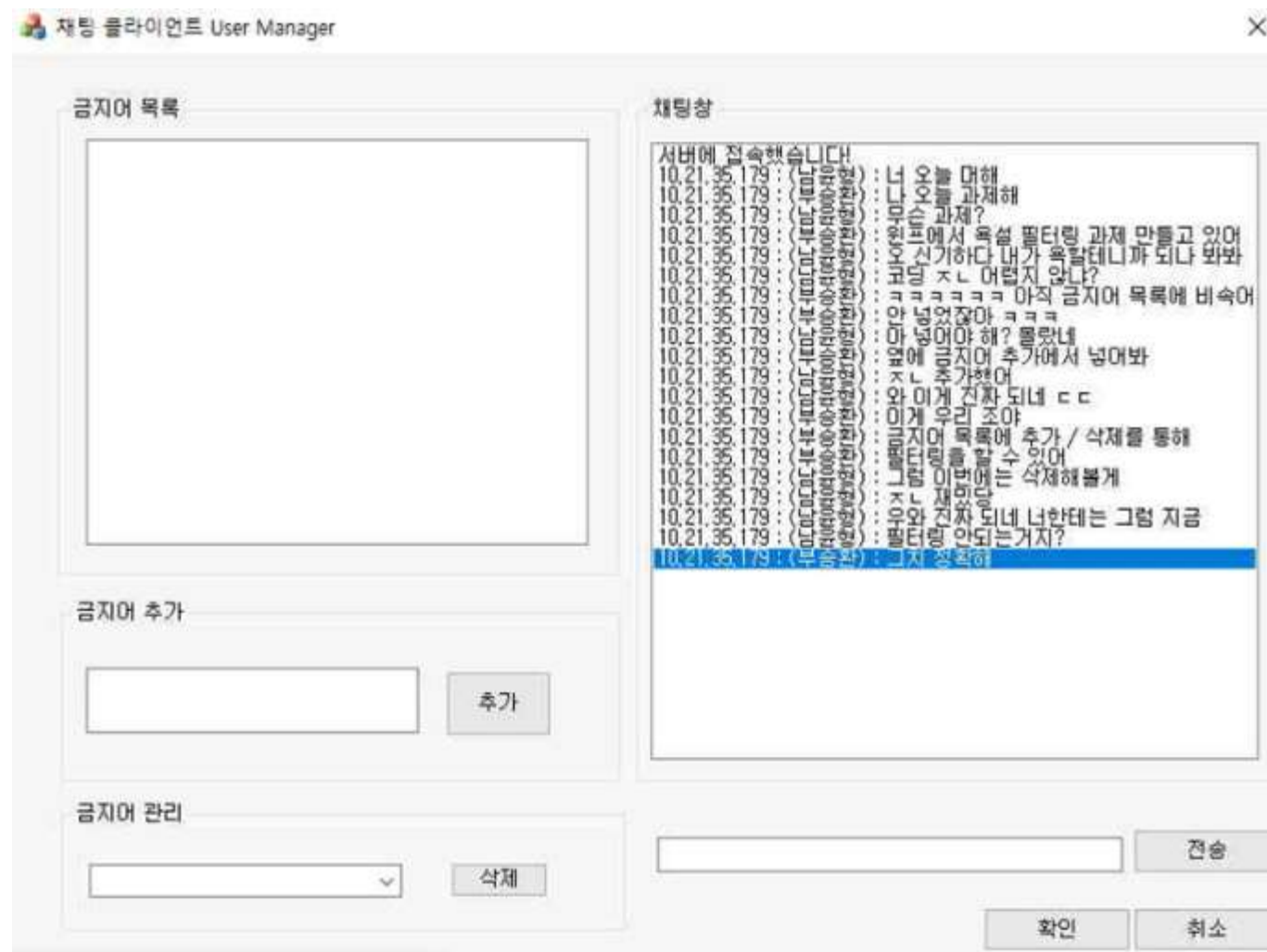


실험 결과 분석

금지어 목록에 추가한 경우 제대로 필터링이 된 모습을 확인할 수 있으며 금지어 목록에서 제외할 경우 원래는 필터링되면 단어가 다시 필터링이 되지 않는 모습을 확인할 수 있다.

분석 결과 - 1

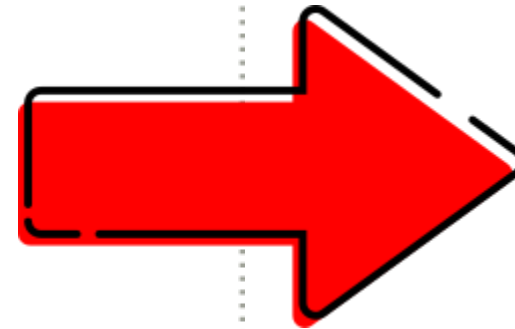
분석 결과 - 2



결론

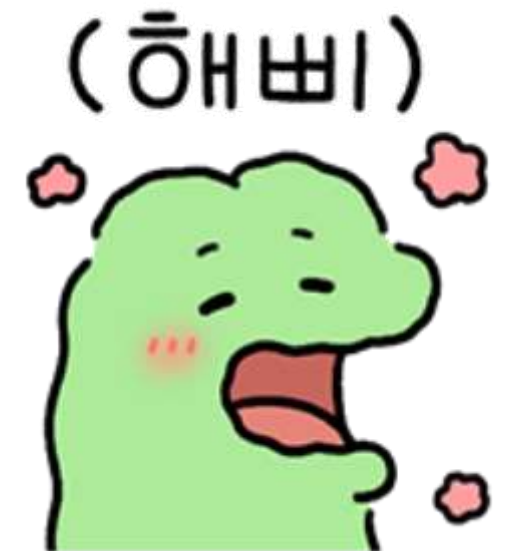
결론 - 1

실시간 채팅 시스템에서 부적절한 언어 사용으로부터 사용자를 보호하고, 건전하고 긍정적인 온라인 커뮤니케이션 환경을 조성하는데 기여할 수 있을 것으로 기대됨

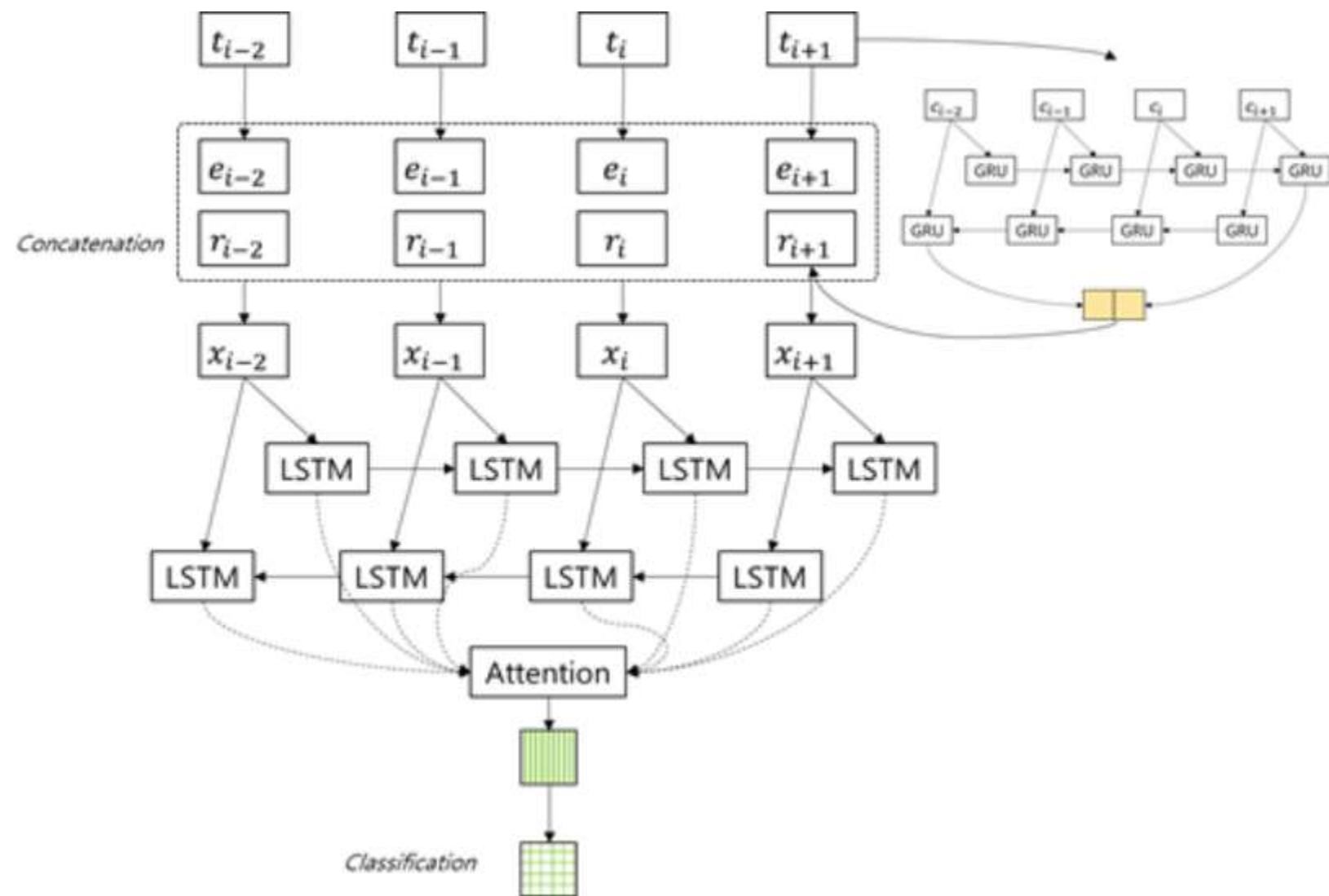


기대효과

이를 통해 온라인 플랫폼 운영자는 윤리적이고 신뢰할 수 있는 서비스를 제공 하는데 도움을 받을 수 있을것으로 보임



향후 연구 방향



- 딥러닝 기반 필터링 시스템을 활용한 복잡한 문맥 인식 시스템
- 인공지능망을 활용한 동적 필터링 기술을 추가적으로 탐구

Q&A

**Thank
you!**