

2025.05.14.(수)

# NS3 Code Generator

: LLM 기반의 네트워크 시뮬레이션 코드 자동 생성

---

202101109 박수화  
202102722 최희원

# 목차

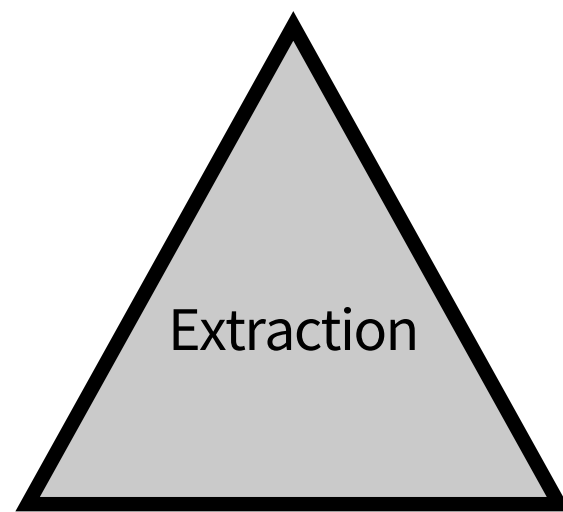
**1 RAG DB 구축 흐름**

**2 접근 방식**

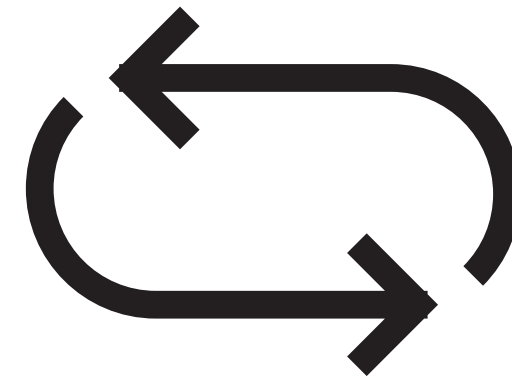
# RAG DB 구축 흐름



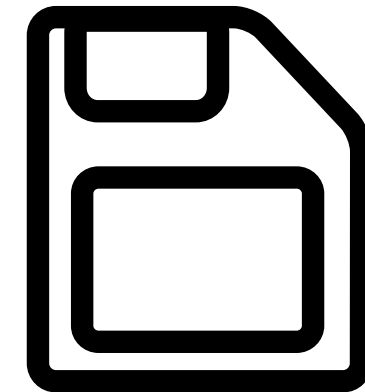
수집



추출



변환



저장

# 접근 방식

## Data 추출 기준

### ① 크롤링은 하되 구조화가 가능만 문서만 골라서 처리

- 공식 문서나 블로그 중에서 헤더나 소제목이 잘 되어있는 글만 처리
- 마크다운 언어 활용해서 분리

### ② 논문 파일(pdf 등) 기반 문서 처리

- 논문, 공식 문서에서 파일 직접 수집
- LLM(KeyBERT 등) 활용해 설정 관련 키워드 추출 → 자동 태깅 또는 웹 검색 쿼리 생성에 활용

# 접근 방식

## RAG DB 구축

### 1. NS-3 에서 자주 등장하는 입력 설정 커버용 키워드 정리

- 논문 5편 다운로드 및 텍스트 추출
- pdf\_keyword.py 를 통해 키워드 추출 + 저장 (json 형식)
- 문제점: 추출된 키워드가 NS-3 코드와 직접 연관이 낮음

### 2. 수동 정의 키워드 기반 자동 쿼리 생성

- 예시 키워드: 802.11ac, UdpEchoClient, PropagationDelayModel 등
- query.py site:medium.com 등과 조합하여 쿼리 텍스트 자동 생성

```
802.11ac site:medium.com
UdpEchoClient site:velog.io
NetDeviceHelper site:tistory.com
```

### 3. 구조화된 문서만 크롤링

→ 예를 들어 <h1> 이런 마크다운 언어 단위로 청킹

### 4. 소제목 같은 것은 자동으로 태그를 통해 매핑

### 5. 임베딩 후 DB에 저장

# 접근 방식

## 1. NS-3 에서 자주 등장하는 입력 설정 커버용 키워드 정리

pdf\_keyword.py

```
1  import os
2  import fitz # PyMuPDF
3  from keybert import KeyBERT
4
5  # PDF 파일들이 있는 폴더 경로
6  pdf_dir = "../rag_db_1/papers/"
7
8  # 모델 초기화
9  kw_model = KeyBERT()
10
11 # PDF 폴더 내 모든 파일 반복
12 def extract_section(text, keywords):
13     lines = text.split("\n")
14     for i, line in enumerate(lines):
15         if any(k.lower() in line.lower() for k in keywords):
16             return "\n".join(lines[i:i+15])
17     return ""
18
19 for filename in os.listdir(pdf_dir):
20     if filename.endswith(".pdf"):
21         pdf_path = os.path.join(pdf_dir, filename)
22         print(f"\n처리 중: {filename}")
23
24         # 1. 텍스트 추출
25         doc = fitz.open(pdf_path)
26         full_text = "\n".join([page.get_text() for page in doc])
27
28         # 2. Simulation 관련 섹션 추출
29         setup_text = extract_section(full_text, ["Simulation", "Experiment", "Parameter", "Setup"]) # 또는 full_text(전체 내용) 사용
30
31         if not setup_text.strip():
32             print("관련 섹션이 없음")
33             continue
34
35         # 3. 키워드 추출
36         keywords = kw_model.extract_keywords(setup_text, top_n=15, stop_words='english')
```

```
38     print("추출된 키워드:")
39     for kw, score in keywords:
40         print(f"    - {kw} ({score:.2f})")
41
42     # 4. 추출된 키워드 저장
43     import json
44     output_dir = "../rag_db_1/keywords/"
45     os.makedirs(output_dir, exist_ok=True)
46
47     output_path = os.path.join(output_dir, filename.replace(".pdf", ".json"))
48     keyword_list = [kw for kw, _ in keywords]
49
50     with open(output_path, "w") as f:
51         json.dump(keyword_list, f, indent=2)
52
53     print(f"저장 완료 → {output_path}")
54
55 # 추출된 키워드 저장
56 import json
57 output_dir = "../rag_db_1/keywords/"
58 os.makedirs(output_dir, exist_ok=True)
59
60 # 저장할 경로 지정
61 output_path = os.path.join(output_dir, filename.replace(".pdf", ".json"))
62 keyword_list = [kw for kw, _ in keywords]
63
64 # JSON으로 저장
65 with open(output_path, "w") as f:
66     json.dump(keyword_list, f, indent=2)
67
68 print(f"저장 완료 → {output_path}")
```

# 접근 방식

## 2. 수동 정의 키워드 기반 자동 쿼리 생성

query.py

```
1  import os
2  import json
3
4  # 1. 경로 설정
5  keyword_dir = "../rag_db_1/keywords/"
6  output_path = "../rag_db_1/keyword_queries.txt"
7
8  # 2. 사용할 site 목록
9  sites = ["medium.com", "velog.io", "tistory.com", "site:blog.naver.com", "site:github.io"]
10
11 # 3. 쿼리 생성
12 queries = []
13
14 for filename in os.listdir(keyword_dir):
15     if filename.endswith(".json"):
16         with open(os.path.join(keyword_dir, filename), "r") as f:
17             keywords = json.load(f)
18
19             for keyword in keywords:
20                 for site in sites:
21                     queries.append(f"{keyword} site:{site}")
22
23 # 4. 결과 저장
24 with open(output_path, "w") as f:
25     f.write("\n".join(queries))
26
27 print(f"쿼리 {len(queries)}개 생성 완료 → {output_path}")
```

# 접근 방식

## 3. 구조화된 문서만 크롤링

```
1  import os
2  import requests
3  from bs4 import BeautifulSoup
4  from langchain.docstore.document import Document
5  from langchain.vectorstores import Chroma
6  from langchain.embeddings import HuggingFaceEmbeddings
7
8  # URL 목록 불러오기
9  with open("../rag_db_1/url_list.txt", "r") as f:
10     urls = [line.strip() for line in f if line.strip()]
11
12  # 크롤링 + 청크 추출 함수
13  def extract_chunks_from_url(url):
14     try:
15         response = requests.get(url, timeout=10)
16         soup = BeautifulSoup(response.text, "html.parser")
17
18         # 구조화된 콘텐츠만 추출
19         content = []
20         for tag in soup.find_all(["h1", "h2", "p"]):
21             text = tag.get_text(strip=True)
22             if len(text) > 50: # 너무 짧은 건 제외
23                 content.append(text)
24
25         # 간단한 청크 나누기
26         chunks = []
27         buffer = ""
28         for paragraph in content:
29             if len(buffer + paragraph) > 500:
30                 chunks.append(buffer.strip())
31                 buffer = paragraph
32             else:
33                 buffer += " " + paragraph
34         if buffer:
35             chunks.append(buffer.strip())
36
37         return chunks
38     except Exception as e:
39         print(f"{url} 에서 오류 발생: {e}")
40     return []
```

builder.py

```
42  # 문서 생성 + 태그 추가
43  documents = []
44  for url in urls:
45     chunks = extract_chunks_from_url(url)
46     for chunk in chunks:
47         doc = Document(
48             page_content=chunk,
49             metadata={"source": url} # 이후 자동 태깅 시 참고
50         )
51         documents.append(doc)
52
53  print(f"총 {len(documents)}개 청크 생성 완료")
54
55  # Chroma DB 저장
56  persist_dir = "../rag_db_1/chroma_db"
57  os.makedirs(persist_dir, exist_ok=True)
58
59  embedding_model = HuggingFaceEmbeddings(model_name="meta-llama/Llama-3.3-70B-Instruct")
60  db = Chroma.from_documents(documents, embedding=embedding_model, persist_directory=persist_dir)
61  db.persist()
62
63  print(f"Chroma DB 저장 완료 → {persist_dir}")
```



# 접근 방식

## 결과

처리 중 : Deep\_Reinforcement\_Learning-Based\_Mobility-Aware\_UAV\_Content\_Caching\_and\_Placement\_in\_Mobile\_Edge\_Networks.pdf

추출된 키워드 :

- mdp (0.46)
- reinforcement (0.40)
- mobile (0.36)
- edge (0.32)
- bandwidth (0.32)
- caching (0.32)
- uav (0.30)
- qoe (0.29)
- algorithms (0.28)
- cache (0.27)
- optimal (0.27)
- learning (0.26)
- unmanned (0.26)
- aerial (0.24)
- network (0.23)

저장 완료 → ../rag\_db/keywords,

처리 중 : 3592149.3592161.pdf

추출된 키워드 :

- simulated (0.43)
- openran (0.40)
- simulation (0.38)
- xapps (0.35)
- 5g (0.34)
- simu (0.33)
- availability (0.31)
- infrastructure (0.31)
- software (0.30)
- platform (0.30)
- distributed (0.28)
- stations (0.28)
- telco (0.25)
- datasets (0.25)
- interface (0.24)

저장 완료 → ../rag\_db/keywords,

처리 중 : electronics-09-00272-v2.pdf

추출된 키워드 :

- simulation (0.38)
- network (0.34)
- fisica (0.20)
- università (0.17)
- ns (0.17)
- bioingegneria (0.16)
- mauro (0.15)
- studi (0.14)
- computer (0.14)
- matematica (0.13)
- systematic (0.13)
- milano (0.13)
- campania (0.12)
- informatica (0.11)
- italy (0.11)

저장 완료 → ../rag\_db/keywords,

처리 중 : 978-3-642-12331-3\_2.pdf

추출된 키워드 :

- networking (0.43)
- protocols (0.43)
- network (0.41)
- simulators (0.40)
- simulator (0.40)
- simulation (0.38)
- ns (0.27)
- testbeds (0.23)
- prototyping (0.23)
- scalability (0.22)
- reproducibility (0.22)
- scalable (0.21)
- reproducible (0.18)
- technologies (0.17)
- research (0.17)

저장 완료 → ../rag\_db/keywords/978-3-642-12331-3\_2.json

- ⇒ 코드 생성에는 도움이 되지 않을 듯한 영양가 없는 키워드가 추출된 느낌
- 논문이 직접적으로 ns3 코드에 사용되는 내용 중심이 아니기 때문 ?

# 접근 방식

## 피드백 및 진행 상황

- 커뮤니티 중심의 코드 수집 (StackOverflow, GitHub, NS-3 커뮤니티 등)
- parameter 설정은 대부분 experiment 섹션 이후에 존재 → 해당 부분 중심으로 수집
- 현재 해당 방향으로 코드 개선 중

# 감사합니다 :)

---

---

---

---