main. py.
@app. on._event ("Start up") ⇒ fastAPI 처음 켜질때 실행.
async def startup_event () :
print ("서버시작~")
initialize_system() → chatbot_logic. py. 의 initialize_system()
print ("시스템 초기화 완료")                    '

Chatbot_logic. Py.
def initialize_system() → 챗봇 시작시 llm, data, vector_db등 필요한 리소스 로딩.

    global llm, retrievers, approved_templetes, rejected_templetes
   → global 키워드를 사용해서 함수 내에서 수정가능하게 함.
    if llm is not None :
        return
   → 중복 초기화 방지
    llm = ChatOpenAI (Model = "gpt-4o", temperature = 0.2)
   → llm 모델 선정, temperature ⇒ 모델의 창의성. ( 0 ~ 1 ), 법규과 가이드 라인을
        (보수적) (창의적)
    지켜야하기 때문에 보수적인 0.2 선택.

    data. dir = ~
    approved_templetes = ~        ⟩ → 참고하는 data 경로 설정.
    rejected_templetes = ~
    from Chromadb.config import settings
    docs_compliance = CustomRuleLoader (os.path.join(data.dir, "compliance_rules.txt")).load()

    class CustomRuleLoader (Base Loader) :
        def __init__ (self file_path : str, encoding : str = 'utf-8') :     Java의 기본생성자와
            self.file_path = file_path                                    같은 기능.
            self.encoding = encoding
        def load (self) → List [Document] :
            docs = []
            try :                                       → data 읽는 형식.        파일을 두라고 변경 지정
                with open (self.file_path, 'r', encoding = self.encoding) as f :
                    content = f.read() → 파일을 utf-8로 읽어서 content에 저장

except File Not Found Error :
    print (f" 경고 ~ )         → 파일 없는경우 예외 처리.
    return [ ]

rule_blocks = re. findall (r'\[규칙시작\] (.?) \[규칙끝\]', content , re. DOTALL)

    re = 정규 표현식 모듈., 문자열에서 특정 패턴 찾기.
    ⌐ r = ` ` 를  RAW string 으로 처리
    ⌐ ∴ txt 에서 '규칙시작', '규칙 끝' 을 찾아서 그 블럭들 리스트 형태로 각각 rule_blocks 에 저장

for block in rule_blocks
    lines = block . strip() . split ('/n')        예) [규칙시작]
    metadata = {}                Source : ∼
    page_content = " "            part : ∼

for line  IS_content_section = False       section. ∼
in lines  If  line.lower(). starts with ('content') :  rule_Id : ∼
        IS_content_section = True       content :
        page_content += line [len('content:' :] . strip t' \n"  핵심정의 ∼
        continue               [규칙 끝]
    if IS_content_section = True :
        page_content += line.strip() + "\n"
    else :
        if ':' in line :
            key, value = line. split (':', 1),
            metadata [key.strip()] = value . strip()
  If page_content :
    docs.append (Document (page_content = page_content . strip(), metadata = metadata))

return docs →  ∴ docs 는  page_content 에 compliant_txt 의 content를 반환
             metadata 에  나머지 데이터를  key : value 값들 저장.

```
docs_generation = CustomRuleLoader( ~ )  → rule.txt를 읽어서
                                          page_content = " ~ "
                                          metadata = { }     를 반환
```

```
docs_whitelist = [ Document (page_content = t) for t in approved_templates]
docs_rejected  = [ Document (page_content = t) for t in rejected_templates]
```

Document 객체는 LangChain에서 텍스트 데이터를 관리하는 표준형식.

풀어서 작성해보면.
```
docs_rejected = []
for t in approved_templates:
    new_doc = Document (page_content = t)
    docs_whitelist.append (new_doc)
```

여기서 approved_templates = load_by_line( ~ ). → 한 줄 씩 갈라서 읽기.
rejected_templates = load_by_seperator( ~ ) '__'를 기준으로 나눠서 읽기.

```
embeddings = openAIEmbeddings (model = ~ ).
Vector_db_path = ~
Client_settings = Settings (anonymized_telemetry = False)
```
→ 뭔 준비

↳ 익명화된 사용 현황 데이터. Chroma_db 기본설정에
  보내지 않기, 사용데이터가 자체 내부의 데이터 이기
  때문에.

```
def create_db (name, docs)
    if docs:
        return Chroma.from_documents (docs, embeddings, ~ )  → 벡터 DB 생성
    ⇒ return Chroma (collection_name = name, embedding_function = embeddings, ~ ).
```
↳ 불명확한 코드, data에 없지만 Chroma_db에만 있는 데이터가 있으면 사용됨.

```
db_compliance = ~
db_ ~
db_ ~
db_rejected = ~
```
→ 벡터 DB

```
def create_hybrid_retriever (vectorstore, docs):
    if not docs:
        return vectorstore.as_retriever (search_kwargs = {"k" = 6})
```
docs가 없을 경우 대비
바이패스.

Vector_retriever = VectorStore.as_retriever (Search_kwargs={"k"=5)
→ 의미기반 검색. 문맥이 비슷한 문서 5개 찾음. "k"는 값을 바꿔면서 최적의 값 찾기.

keyword_retriever = BM25Retriever.from_documents (docs)
→ 키워드 기반 검색. 요청에 포함된 키워드와 들어있는 문서 찾기.

keyword_retriever.k = 5 → 상위 5개. 즉 변경 가능

Ensemble_retriever = EnsembleRetriever (retrievers [vector_retriever, keyword_retriever],
                        weights = [0.5, 0.5])

→ 두 검색기를 하나로 합쳐서 순위 매김, weights = [0.5, 0.5]는 가중치 50.50을
   주어서 순위 부여.

if Ranker :
        compressor = FlashRank Rerank (top_n=3).
        return ContextualCompression Retriever ( base_compressor = compressor, base_retriever=
                        ensemble_retriever)

return  ensemble_retriever.
→ 재정렬 : ensemble_retriever 에서 부여된 순위를 다시 검토. ensemble_retriever에
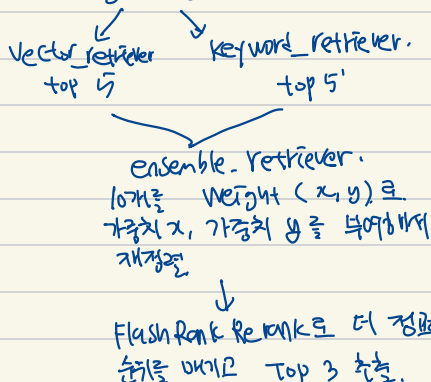    현재 10개가 들어 있고 거기서 top 3를 선정해 그것만 추출.

retrievers ['compliance'] = Create_hybrid_retriever (db.compliance., docs_compliance)
  :   ['generation'] :          =
  :   ['whitelist'] =            :
  :   ('rejected') =             =
→ 4개의 검색기를 생성. 과정은 : 사용자 요청

                    Vector_retriever        keyword_retriever.
                    top 5                   top 5'

                            ensemble_retriever.
                            10개를 weight (x, y)로
                            가중치 x, 가중치 y를 부여해서
                            재정렬
                                ↓
                            FlashRank Rerank로 더 정밀하게
                            순위를 매기고 Top 3 추출.

math. 1차로 복기. 서버 실행, 초기 리소드 로딩 완료.

채팅 시작.

프론트에서    /api/chat 으로  post. 요청.

```
@ app.post ("/api/chat")
async def chat (request : ChatRequest) -> Dict:    ) → 채팅 시작.
    session_state = request.state.
```

```
Class ChatRequest (BaseModel):          } 중요성 예시.
    Message : Str
    State : Optional [Dict] = Field (default_factory =dict)
```

Message ⇒ 필수 입력값.

State ⇒ 선택 입력, 만약 입력 → Dict 형식으로 만듦, 만약 !입력 → 빈 Dict 사용

```
If not session_state:          세션이 비어있다면 초기화.
    session_state = {
        'step' =
            :
        'correction_attempts' : 0 }
response_data = process_chat_message (request.message, session_state) → 채팅 메시지 처리로직.
```

    chat_logic.py.

```
def process_chat_message (message : str, state : dict) → dict:

    If   State ['step'] == 'Initial':
        State ['original_request'] = message
        State ['step'] = 'recommend_templates'
        similar_docs = retrievers ['whitelist'].invoke (message).
            ↳ p.4의 검색기 참고. 실행버튼    ↳ 사용자의 첫 요청.

        if not similar_docs :
            State ['step'] = 'select_style'
            return {
                'message' : '유사한 기존 템플릿~'
                'state' :
                'options' : ['기본형', '이미지형', '아이템리스트형']
            } → 유사템플릿 못찾으면  이후단계 스킵.
```

→ 유사 템플릿 존재

```
templates = [f' 템플릿 {i+1} : \n {doc.page_content}" for i, doc in
                enumerate (similar_docs [:3])]
```

⇒ for i, doc in enumerate (similar_docs [:3]):
    f " 템플릿 {i+1} : \n {doc.page_content}"

return {
        'message' : ~
        'state' : state
        'options' : [ '템플릿1', '템플릿2', '템플릿3', '신규생성' ]
        'templates' :
    }

⇒ state 에는 변경된 step, original_request 가 당되서 넘어감.
              recommend_template,   ↳사용자 요청
                                    (추가)

```
elif state ['step'] == 'recommend_templates' :
    if message in [ '템플릿1', '템플릿2', '템플릿3' ] :
        template_idx = int (message.split() [1]) - 1
        similar_docs = retriever ['whitelist'] . invoke (state['original_request'])
```
⇒ similar_docs 를 이미 추출했는데 또 검색기로 추출하는 이유.
    BM25, text_embedding, 백터검색, 순위정렬 (ensemble, flash rank)
    모두 정해진 계산 방식에 의해 검색되고 정렬되므로 결과가 달라지지 않음.
    따라서 템플릿을 모두 다 저장하는 것보다 original_request 만 저장하고
    사용시에만 다시 실행.

```
    state ['template_draft'] = similar_docs [template_idx]. page_content.
```

```
    elif message == '신규생성'
        state ['step'] = 'select_style'
        return {
            'message' : '새로운 ~ '
            'state' : state
            'options' : [ '기본형', '이미지형', '이미지리스트형' ]
    else. state ['step'] = 'select_style'
        return process_chat_message (message, state)
```

```python
if  State.get ('step') = 'select _style' :
      if  message. in  ['기본형', '이미지형', '아이템리스트형'] :
           State ['selected _style'] = message.
      else :
           State ['selected _style'] = '기본형'


      State ['step'] = 'generate_ and _validate'
      return  process _ chat _ message ( message, state)



if  State .get ('step') == 'generate _and _validate' :
      template _draft = generate _template ( state ['original _request'], state.get ('selected_
                                        style', '기본형'))
 def  generate _template ( request : str, style : str) → str :
        example _docs = retrievers ('writelist') . invoke (request)
        examples ="\n\n" .join ([f"예시 {i+1} : \n { doc. page _content}" for i, doc in
                                            enumerate (example _docs)])
        expansion _prompt = Chat Prompt Template. from_template (
              """"  ~  """" ) ⇒ 템플릿 초안 작성 명령.
        expansion _chain = expansion _prompt ) llm | Str Output Parser()
        expanded _draft = expansion _chain . invoke (
              {                                                      → 작업 명령서
                  "original _request" : request,                     → 작성자
                  "style" : style,                                   → 최종 점수 and 표장
                  "examples" : examples
              })
        return expanded _draft.
      State ['template _draft'] = template _draft
      validation _result = validate _template (template _draft)
```

```python
def validation_template (draft : Str) → dict :
    parser = JsonOutput ~  =) 출력 형식 준비
    Step_back. prompt = Chat prompt Template. from _template (.
                "''''" ~ "'''")
    Step_back _question = Step_back. prompt | llm ( StrOut put parser), Invoke (9 "draft"
                                                                    : draft})

    compliance - docs = retrievers ['compliance'] . Invoke (step_back. question)
    rules _with _metadata ="\n\n" .join ( ~ )  =) ~ 검색기에서 출력한 3개의
                                                    규칙 list로 저장

    rejected _docs = retrievers ['rejected']. Invoke (draft)
    Rejections = " \n\n" . join (~ ) =) list로 저장
    validation _prompt = Chat prompt Template . from _template (
            "''''"  ~ "'''" )
    validation_chain = validation _prompt | llm | parser
    result = validation _chain .Invoke ( ~ )
    return result


State ['correction _attempts'] =0
If validation _result ['status'] == "accepted" :
    State ['step'] = "completed"
    return process_chat _message (message, state) => step이 completed로
                                                    내가이어서 다음 단계로 이동
else
    State ['step'] = 'correction' =) step이 correction으로 변경.
    return {                        수정단계로 이동.
        'message' : f' 템플릿 ~ '
        'state' : State
    }
elif State ['step']== 'correction' :
    If State [' correction _attempts'] < MAX _CORRECTION-ATTEMPTS :
        corrected _template = Correct _template (state)
```

```
def    correct _template (State : dict ) → Str :
       attempts = State. get (' correction _attempts', 0) ⇒ correction. attempts
                                                        개체이. 값으면 기본값0.
         If attempts == 0 :
              Instruction = "3. ~  "
         elif  attempts == 1 :
              Instruction = " 3, ~  "
         else :
              Instruction = """"  3. ~    ""  "
       correction _prompt _template = ""  "   ~   "" "
       correction _prompt = ChatPromptTemplate. from _template ( correction_prompt
                                                                      _template)

       correction _prompt = Correction _prompt . partial ( dynamic _instruction =
                                                                      instruction )


       correction _chain = correction _ prompt | llm | StrOutput Parser ()
       new _draft = correction _chain. Invoke ({" original _ request" : state [
                                              ' original_request' ] ~ })

         return  new _draft.
```

시도 횟수별  프롬프트 변화
→ 광고성문구 제거
   ↓
   쿠폰, 할인 등과 같은 키워드제거
   ↓
   정책에 맞게 과감한 변화

```
State [' template _draft'] = corrected - template.
State [' correction _attempts'] += 1
validation _result = validate _template (corrected _template) ⇒ 대재검증
State [' validation _result' ] = validation _result
If  validation _result [' status' ] == "accepted":
       State [' step'] = " completed"
       return  process _chat. message (message, state) ⇒ completed
else :                                                        단계의 응
       return   process _chat _message (message, state) ⇒ correction
                                                           단계로
                                                         재계 호출
else: ← MAX _CORRECTION _ATTEMPTS 넘으면 사용자에게
                                  수동 수정 요청.
```

```
          state ['step'] = 'manual_correction'
          return {
                   'message' : ~
                   'state' : ~
                   'options' : ['포기하기']
                }
elif   state ['step'] = 'manual_correction':
       if  message == '포기하기':           포기하면 초기단계로 이동
            state ['step'] = 'Initial'
            return {
                    'message' : ~
                    'state' : ~
                 }
       else :
             state ['template_draft'] = message
             validate_result = validate_template (message)
                             검증단지
             state ['validation_result'] = validation_result
             if  validation_result ['state'] == 'accepted' :
                  state ['step'] = "completed"
                  return   process_chat_message (message, state)
             else :
                  return {
                          'message' : ~
                          'state' : ~
                          'options' : ~
                       }
elif   state ['step'] == 'completed' :
       final_template = state.get ("template_draft", "")
       html_preview = render_final_template (final_template)
```

```
def render_final_template (template_string : str) → Str :
```
}
→ 일정한 형식의 html로 변경

```
    return html_output
```

```
    parameterized_result = parameterized_result (final_template)
```

템플릿에서 변수로 선택 가능한 부분은 # 변수명으로 선정.

```
    return {
            "message" :
             "State" :
             'template' :
             ' html_preview'
             ' editable_variables'
          }
```

```
return {
        'message' : ~
         'State' : ~
      }
```
) process_chat_message 종화.