



LLM을 이용한 COLREGs 기반 자율운항 의사결정 알고리즘 개발

서진혁¹·임승현²·정성현¹·임선혁¹·김시원¹·김연수¹·김정현¹·박종용^{1,3,†}
국립부경대학교 마린융합디자인공학과 조선해양전공¹
한국해양과학기술원 부설 선박해양플랜트연구소²
국립부경대학교 조선해양시스템공학과³

Development of an LLM-based Decision-making Algorithm for COLREGs-compliant Autonomous Navigation

Jin-Hyeok Seo¹·Seung-hyeon Lim²·Seong-Hyeon Jeong¹·Sun-Hyuck Im¹·Si-Won Kim¹·Yeon-Soo Kim¹·
Jeong-Hyeon Kim¹·Jong-Yong Park^{1,3,†}
Department of Marine Design Convergence Engineering Pukyong National University¹
Korea Research Institute of Ships & Ocean Engineering²
Department of Naval Architecture and Marine System Engineering Pukyong National University³

This is an Open-Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License(<http://creativecommons.org/licenses/by-nc/3.0>) which permits unrestricted non-commercial use, distribution, and reproduction in any medium, provided the original work is properly cited.

Interpreting qualitative terms in the International Regulations for Preventing Collisions at Sea (COLREGs) remains challenging for conventional geometric collision-avoidance algorithms, while large language models (LLMs) can exhibit geometric hallucinations that lead to unsafe maneuvers. To address this, we propose a hybrid decision-making framework that couples LLM-based semantic reasoning with a deterministic Safety Shield. The framework employs Retrieval-Augmented Generation (RAG) to ground decisions in regulatory texts and authoritative interpretations, while the Safety Shield enforces hard physical constraints based on the Closest Point of Approach (CPA) and Time to CPA (TCPA) to ensure collision-free navigation. Validation on real-world Automatic Identification System (AIS) encounters shows that the proposed method achieves the lowest Intrinsic Illegality Rate (IIR) among baselines, remaining robust under sampling uncertainty. Qualitative case studies further demonstrate compliant behavior in complex scenarios, including speed-controlled pass-astern maneuvers and unified actions that resolve multi-vessel conflicts. By structurally decoupling semantic interpretation from physical safety assurance, this study provides a robust pathway toward explainable and COLREGs-compliant autonomous navigation.

Keywords : Large language models(거대언어모델), Retrieval-augmented generation(검색증강생성), International regulations for preventing collisions at sea(국제해상충돌예방규칙), Explainable artificial intelligence(설명 가능한 인공지능)

1. 서론

해상 운송은 전 세계 교역량의 80% 이상을 담당하는 글로벌 공급망의 핵심 인프라이다. 그럼에도 해양 사고는 지속적으로 발생하며, 사고 원인의 상당 비중이 인적 과실(human error)과 관련 된다는 보고는 해사 산업이 여전히 인간 의존적 위험에 취약함을 시사한다(Allianz Global Corporate & Specialty (Allianz), 2025). 이러한 위험을 구조적으로 완화하기 위한 대안으로 자율

운항선박(Maritime Autonomous Surface Ships, MASS) 기술이 대두되고 있다. 그러나 실제 해역 도입을 위해서는 단순한 경로 생성만으로는 충분하지 않으며, 복잡한 조우 상황을 해석하고 적절한 조치를 결정·검증하는 의사결정 능력이 필수적이다. 특히 국제해상충돌예방규칙(International Regulations for Preventing Collisions at Sea, COLREGs)은 “선박 운항술(good seamanship)”과 같은 정성적·맥락 의존적 표현을 다수 포함하고 있어, 기하학적 수치뿐 아니라 상황의 ‘맥락(context)’을

반영한 유연한 판단을 요구한다(IMO, 1972).

기존의 속도 장애물(Velocity Obstacle, VO) (Fiorini and Shiller, 1998) 또는 모델 예측 제어(Model Predictive Control, MPC) 기반의 결정론적 접근은 재현 가능하고 검증 가능한 출력을 제공한다는 장점이 있으나, COLREGs 준수 행동 선택을 명시적으로 고려하는 경우에도(Johansen et al., 2016) 다선박·복합 조우 상황에서 정성 규정의 적용 근거를 일관되게 설명하기 어렵고, 안전 여유를 크게 설정하는 방식으로 귀결되어 효율과 안전 사이의 상충 관계(trade-off)를 효과적으로 조율하지 못하는 한계가 지적되어 왔다(Vagale et al., 2021). 최근에는 이를 보완하기 위해 거대언어모델(Large Language Models, LLMs)을 활용하여 조우 상황을 자연어로 해석하고 규정 적용을 시도하는 연구가 증가하고 있으나(Zhang et al., 2023), LLM의 환각(hallucination) 및 기하학적 추론 오류는 물리적 안전이 전제되어야 하는 선박 운항에서 치명적 위험요소가 될 수 있다(Ji et al., 2023).

본 연구는 이러한 딜레마를 완화하기 위해, LLM의 유연한 맥락 추론 능력과 결정론적 안전 통제를 결합한 하이브리드 의사결정 프레임워크를 제안한다. 특히 모델의 판단 신뢰성을 높이기 위해 검색 증강 생성(Retrieval-Augmented Generation, RAG) 기법을 적용하였다(Lewis et al., 2020). RAG란 모델이 자체 지식에만 의존하지 않고 외부의 신뢰할 수 있는 지식 베이스를 실시간으로 참조하여 답변을 생성하는 기술로, 환각 현상을 효과적으로 완화한다(Ji et al., 2023). 본 프레임워크는 이를 통해 COLREGs 원문과 공신력 있는 해설을 실시간으로 검색하여 참조하게 함으로써, 모델의 자의적 해석에 따른 환각 현상을 배제하고 철저히 규정 텍스트에 기반한 항해 판단을 유도한다(Lewis et al., 2020). 또한 LLM이 제시한 회피 조치가 최근접 거리(Closest Point of Approach, CPA) 등 물리적 안전 제약을 만족하는지 결정론적 안전 쉴드(safety shield)로 이중 검증하여, 언어적 추론 결과가 기하학적 안전 조건을 위반하는 경우를 체계적으로 배제한다. 즉, 규칙의 ‘해석’은 LLM이 수행되되 물리적 ‘안전’

은 안전 쉴드가 강제하는 역할 분리를 통해 설명 가능성과 안전성을 동시에 확보한다.

본 연구는 COLREGs 원문과 공신력 있는 해설을 대상으로 RAG 기반 근거 검색을 수행하여, 판단이 규정 텍스트에 정렬되도록 구성하였다. 또한 LLM 출력에 행동, 적용 규칙, 조우 해석, 역할, 근거 인용, 설명을 포함하는 구조화 스키마를 강제함으로써 판단 과정의 추적 가능성을 확보하였다. 나아가 CPA/TCPA 기반 Safety Shield를 런타임 안전 계층으로 결합하여, 언어모델의 기하학적 추론 오류가 물리적 안전 위반으로 전이되는 실패 모드를 차단하였다.

2. 제안 방법

2.1 프레임워크 구성

본 절에서는 조우 상황에서 COLREGs 판단을 근거 인용이 가능한 설명 형태로 생성하기 위한 입력-출력 정의와 전체 처리 흐름을 기술한다. 제안 프레임워크는 (i) 규정 근거 검색(RAG), (ii) 근거 인용 기반 판단 생성(LLM), (iii) 결정론적 안전 쉴드로 구성되며, Fig. 1에 전체 아키텍처를 제시한다.

시각 t 에서 자선을 o , 주변 선박 집합을 N_t 로 두며, 조우 상황 입력 s_t 는 식 (1)과 같다.

$$s_t = (X_t^{(o)}, X_t^{(i)} | i \in N_t, \varepsilon_t) \tag{1}$$

여기서 $X_t^{(o)}$ 는 자선의 위치·침로·속력, $X_t^{(i)}$ 는 타선 i 의 위치·침로·속력, ε_t 는 항로 경계나 제한수역 등 상황 해석에 필요한 환경 정보를 의미한다. 본 연구에서는 AIS로부터 획득한 위치(위·경도), SOG, COG를 사용하여 주변 선박의 운동 상태를

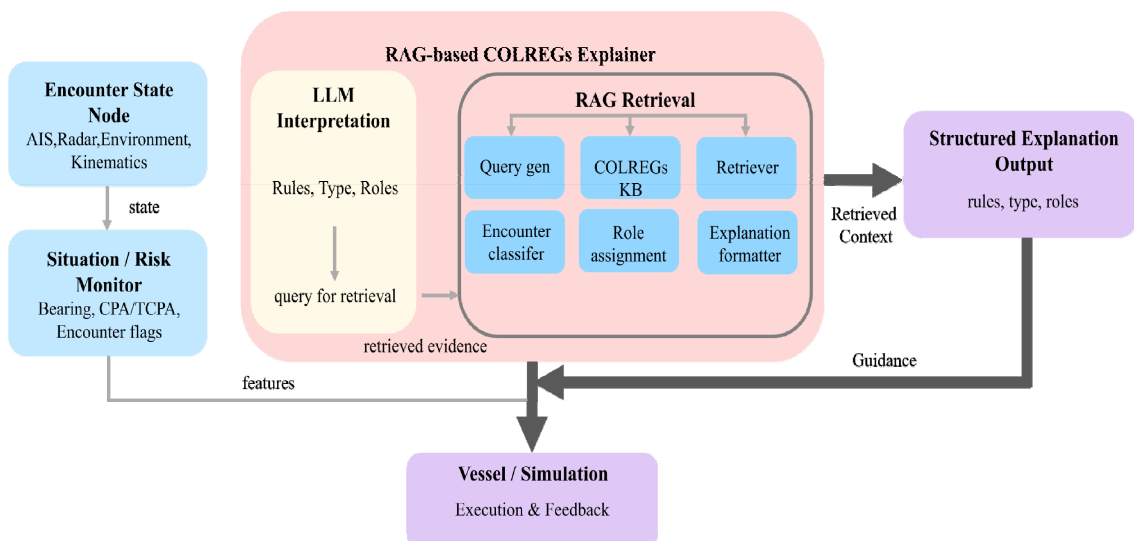


Fig. 1 Overall architecture of the proposed RAG-based COLREGs explainer

구성한다. 이때 LLM 입력 집합 N_t 는 모든 주변 선박을 포함하지 않고, 매 시점 t 에서 자선(ego)과 주변 선박 i 간 최근접 접근 거리(DCPA) 및 최근접 접근 시간(TCPA)을 산출한 뒤, $TCPA_i(t) > 0$ 인 선박 중 DCPA 오름차순으로 정렬한 상위 k 척만을 동적으로 필터링하여 구성한다.

출력 y_t 는 LLM이 제안한 행동과 그 판단 근거를 구조화한 결과로 정의하며, 식 (2)와 같다.

$$y_t = (a_{llm}, L_t, k_t, p_t, C_t, r_t) \quad (2)$$

식 (2)에서 a_{llm} 는 LLM이 제안한 행동(action), L_t 는 적용 규칙(조항 번호/규칙명) 집합(rules), k_t 는 조우 유형 해석(encounter), p_t 는 유지선/피항선 관점의 역할 해석(role), C_t 는 RAG로 검색된 규정 문장 중 인용된 근거(citations), r_t 는 자연어 설명(rationale) 텍스트를 의미한다. Fig. 1은 입력 s_t 로부터 규정 근거 검색과 판단 근거 생성이 수행되는 전체 처리 흐름을 나타낸다. 프레임워크의 핵심 요소인 COLREGs KB(Knowledge Base)는 국제해상충돌예방규칙의 전문(full text)과 관련 해설서(commentaries)를 텍스트 청크(chunk) 단위로 구조화하여 구축된 데이터 저장소이다. 시스템은 이 지식 베이스를 활용하여 현재 조우 상황에 부합하는 규칙 원문을 실시간으로 검색하고 프롬프트에 주입함으로써, LLM의 추론을 실제 법적 텍스트에 정렬(grounding)시키는 ground truth 역할을 수행한다.

2.2 검색 증강 생성(RAG)

본 연구는 COLREGs 원문과 공신력 있는 해설·지침을 문서 집합 D 로 구성하고, 이를 지식베이스로 사용한다. 검색 단계는 현재 조우 상황과 연관된 규정 문장을 후보로 제한함으로써, 생성 과정이 규정 텍스트에 기반하도록 만드는 역할을 수행한다.

조우 상태 s_t 로부터 질의 $q_t = q(s_t)$ 를 구성하고, 지식베이스 D 에서 관련성이 높은 근거 후보 집합 Z_t 를 검색하는 과정은 식 (3)과 같다.

$$Z_t = R(D, q_t) \quad (3)$$

식 (3)에서 $R(\cdot)$ 는 근거 코퍼스 D (COLREGs 조항/해설 등)의 문서들을 사전에 임베딩해 두고, 매 시점 생성되는 질의 q_t 를 임베딩 벡터로 변환한 뒤 코사인 유사도(cosine similarity)를 기준으로 상위 K 개의 근거 후보 $Z_t = \{z_t^1, \dots, z_t^K\}$ 를 반환하는 연산자이다. 각 k_t^j 는 Evidence ID, 출처(문서/조항), 그리고 핵심 스니펫으로 구성된다.

2.3 근거 인용 기반 판단 설명 생성(LLM)

생성 단계에서는 조우 상황 s_t 와 근거 후보 Z_t 를 입력으로

받아 출력 y_t 를 생성하며, 식 (4)와 같이 표현된다.

$$y_t = G_\theta(s_t, Z_t) \quad (4)$$

여기서 $G_\theta = (\cdot)$ 는 파라미터 θ 를 갖는 LLM 기반 생성 함수로, 입력(s_t, Z_t)를 “상황 요약-근거(evidence)-출력 스키마” 형태의 프롬프트로 구성한 뒤 JSON 출력 y_t 를 생성한다. 즉, $G_\theta : (s_t, Z_t) \rightarrow y_t$ 로 정의되며, 생성된 출력은 후처리를 위해 구조화된 필드로 파싱된다. 출력은 후처리 및 평가를 위해 JSON 객체로 강제하며, 필수 필드는 action(a_{llm}), rules(L_t), encounter(k_t), role(p_t), citations(C_t), rationale(r_t)로 고정한다. citations에는 근거 후보 Z_t 의 식별자(evidence ID)와 출처를 명시하여, 모델의 규칙 선택과 조치 서술이 입력된 근거에 의해 지지되는지 추적 가능하도록 구성한다.

프롬프트는 “상황 요약-근거-출력 스키마”의 순서로 구성한다. 상황 요약에는 조우 유형, 유지선/피항선 역할, 상대선 방향, 현재 분리거리, 예측 최소 분리거리 등 판단에 필요한 최소 정보를 포함한다. 근거 Z_t 는 evidence 1, evidence 2와 같이 번호를 부여하고 문서명과 스니펫을 함께 제공한다.

출력 정렬을 위해 모델은 Z_t 에 포함된 근거만 인용하도록 제한하며, 충분한 근거가 검색되지 않은 경우에는 citations를 비우고 rationale에 근거 부족을 명시하도록 유도한다. 만약 JSON 파싱에 실패하거나 필수 필드가 누락될 경우 1회 재출력을 수행하며, 재시도 후에도 실패하면 보수적 기본 출력을 사용하되 원시 출력은 로그로 저장한다. 생성된 행동은 2.5절의 안전 실드를 통해 최종 검증된다.

2.4 LoRA(+SFT)

본 연구의 기본 모델(base model)로는 오픈 소스 기반의 고성능 언어 모델인 LLaMA(Touvron et al., 2023)를 채택하였다. LLaMA는 다양한 벤치마크에서 입증된 우수한 추론 능력을 바탕으로 복잡한 항해 상황을 해석하는 데 적합한 기초 성능을 제공한다. 그러나 대규모 파라미터를 가진 모델 전체를 미세조정(Full Fine-Tuning)하는 것은 막대한 연산 자원을 필요로 한다.

이에 본 연구는 LoRA(Hu et al., 2021) 기법을 적용하여 효율적인 학습을 수행하였다. LoRA는 사전 학습된 모델의 가중치는 고정된 채, 학습 가능한 저랭크 행렬(low-rank matrices)만을 주입하여 업데이트하는 방식이다. 이를 통해 NVIDIA RTX 5090과 같은 단일 GPU 환경에서도 메모리 요구량을 최소화하며 효과적인 지도 미세조정(Supervised Fine-Tuning, SFT)이 가능하도록 최적화하였다.

학습 데이터는 ‘조우 상황 입력’과 ‘전문가 정답(expert label)’ 쌍으로 구성된다. 여기서 전문가 정답은 인간 판단의 주관성과 불일치를 배제하기 위해, 검증된 규칙 기반 알고리즘(rule-based algorithm)의 결정론적 출력을 사용하여 생성하였다. 즉, 본 학습

은 새로운 규칙 해석 논리를 생성하는 것이 아니라, 규칙 기반 알고리즘의 판단을 언어 모델이 일관되게 재현하도록 정렬 (alignment)하는 것을 목적으로 한다.

구체적인 학습의 핵심 목표는 (i) 정의된 JSON 스키마의 엄격한 준수(format stabilization), (ii) 제공된 법규 텍스트 범위 내에서의 근거 인용을 통한 환각(hallucination) 억제, 그리고 (iii) 샘플링 불확실성 하에서도 일관된 행동과 규칙을 선택하는 출력 정렬이다. 이러한 최적화를 위해 사용된 학습 하이퍼파라미터는 Table 1에 요약하였다. 구체적으로, LoRA의 어댑터 크기를 결정하는 rank(r)와 스케일링 계수(α)는 파라미터 효율성과 학습 안정성을 고려하여 각각 16과 32로 설정하였다. 또한, 배치 크기(batch size)는 제한된 GPU 메모리를 고려하여 8로 설정하되, 그라디언트 누적(gradient accumulation) 기법을 병행하여 안정적인 수렴을 유도하였다.

Table 1 Hyperparameters for LoRA-based adapter tuning (SFT)

Hyperparameter	Setting
hyperploRA (r, α, p)	$r = 32, \alpha = 64, p = 0.1$
Learning rate	2×10^{-4}
batch size	$8 (B = 1, G = 8)$
Epochs	$E = 50$
Optimizer	Adam
Sequence length	1024/512

2.5 안전 쉴드(Safety Shield)

LLM은 COLREGs 해석에는 강하지만, 상대 기하에 기반한 충돌 위험을 과소평가하는 오류(geometric hallucination)가 발생할 수 있다. 이를 보완하기 위해 본 연구는 LLM이 제안한 행동 a_{llm} 을 최종 반영하기 전에 결정론적 안전 쉴드로 위험 여부를 판정하고, 필요 시 안전 행동으로 대체한다.

안전 쉴드는 CPA와 TCPA를 이용해 위험 상태를 판단한다. 시각 t 에서의 조우 상태 X_t 에 대한 위험 상태 함수 $f_{risk}(X_t)$ 는 식 (5)와 같다.

$$f_{risk}(X_t) = \begin{cases} 1, & d_{CPA} < d_{safe} \\ & \text{and } 0 < t_{CPA} < t_{safe} \\ 0, & \text{otherwise} \end{cases} \quad (5)$$

여기서 X_t 는 시각 t 에서의 자선-타선박 상대 상태를 의미하며, d_{CPA} 와 t_{CPA} 는 해당 상태에서부터 계산되는 최근접 접근 거리와 최근접점 도달 시간 이다. 또한 d_{safe} 와 t_{safe} 는 각각 최소 안전거리와 안전 여유시간이며, 본 연구에서는 $d_{safe} = 0.5NM$, $t_{safe} = 6min$ 을 사용하였다.

여기서 a_{shield} 는 사전에 정의된 회피 행동 후보 중, 예측 지평선

내 DCPA/TCPA(또는 예측 최소 분리거리)가 가장 안전한 방향으로 개선되도록 선택된다. 다선박 조우의 경우 모든 타선에 대해 동일한 평가를 수행하고, 가장 불리한 타선 기준(worst-case)에서 안전 여유가 최대가 되도록 a_{shield} 로 결정된다. 이를 위해 본 연구에서는 쉴드 개입 여부를 나타내는 이진 스위치 변수 σ_t 를 도입하였다. σ_t 는 식 (5)에서 산출되는 위험 판정값 $f_{risk}(X_t)$ 과 LLM 제안 행동의 적절성을 참조 하여 매 시점 결정된다. 구체적으로, $f_{risk}(X_t) = 1$ 인 위험 상황에서 LLM 행동이 소극적(passive)일 때 $\sigma_t = 1$ 로 설정하여 a_{shield} 를 적용하고, 그 외에는 $\sigma_t = 0$ 으로 a_{llm} 을 적용한다.

3. 시험 설정 및 평가 방법

3.1 시험 구성

본 연구는 시뮬레이션이 아닌 실제 AIS 로그에서 관측된 조우 상황을 대상으로 오프라인 검증을 수행하였다. 검증 데이터는 선박해양플랜트연구소(Korea Research Institute of Ships and Ocean Engineering, KRISO)가 운영하는 해상시험 실증선 해양누리호(Haeyang Nuri)에서 수집된 AIS 로그로 구성된다. 해상시험 장면은 Fig. 2에 제시하였으며, 해양누리호의 주요 제원은 Table 2에 정리하였다.

각 시각 t 에서의 조우 입력은 2.1절에서 정의한 조우 상태 s_t

Table 2 The principal particulars of Haeyang Nuri.

Item	Value
Length	26.5 m
Beam	5.4 m
Draft	1.4 m
Displacement	97 t
Design speed	12 kn
Main engines	441 kW ×2 (sets)



Fig. 2 Testbed ship Haeyang Nuri

식 (1)로 구성하며, 동일 s_t 에 대해 비교군별 판단 근거 y_t 를 생성하여 성능을 비교하였다. 원시 AIS 시계열은 선박별로 시간 순 정렬 후 1초 간격으로 리샘플링하였으며, 이때 발생하는 AIS 수신 공백으로 인한 누락 데이터는 선형 보간(Linear Interpolation)으로 보정하여 자선(o)과 타선(i)의 조우 상태 s_t 를 구성하였다.

조우 기하 및 예측 최소 분리거리 계산에는 AIS의 위도(latitude), 경도(longitude), SOG, COG를 사용하였다. 이하에서는 위도, 경도, 속력, 침로를 각각 $\phi(t), \lambda(t), u(t), \psi(t)$ 로 표기하며, 자선은 $(\cdot)^{(o)}$, 타선은 $(\cdot)^{(i)}$ 로 구분한다. 위·경도는 degree($^\circ$) 단위로 제공되므로, 근거리 조우 구간에서는 지구 곡률을 무시한 국부 평면 근사(equirectangular approximation)를 적용하여 각도 차이를 해리(NM) 단위의 평면 변위로 변환하였다. 항해 단위에서 1 NM는 위도 1 arcminute에 해당하므로 1° 위도 차이는 60 NM이며, 경도 방향 변위는 위도 ϕ 에서 자오선 수렴에 의해 $60\cos(\phi)$ [NM/deg]로 축척된다. 따라서 동(East)-서(West) 방향 성분 $\Delta x(t)$ 와 북(North)-남(South) 방향 성분 $\Delta y(t)$ 는 식 (6)과 같이 근사된다.

$$\begin{aligned} \Delta x(t) &= 60\cos\left(\frac{\pi}{180}\phi^{(o)}(t)\right) \\ &(\lambda^{(i)}(t) - \lambda^{(o)}(t)), \end{aligned} \tag{6}$$

$$\Delta y(t) = 60(\phi^{(i)}(t) - \phi^{(o)}(t)).$$

속도는 AIS의 SOG/COG로부터 계산하였고, 상대 속도 성분은 식 (7)과 같다. 여기서 $u(t)$ 의 단위는 kn이며, 1 kn은 1 NM/h에 해당한다. 또한 $\psi(t)$ 는 degree($^\circ$)단위이므로 sin, cos 계산을 위해 $\frac{\pi}{180}$ 을 곱해 rad로 변환하였다.

$$\begin{aligned} \Delta v_x(t) &= u^{(i)}(t)\sin\left(\frac{\pi}{180}\psi^{(i)}(t)\right) - \\ &u^{(o)}(t)\sin\left(\frac{\pi}{180}\psi^{(o)}(t)\right), \end{aligned} \tag{7}$$

$$\begin{aligned} \Delta v_y(t) &= u^{(i)}(t)\cos\left(\frac{\pi}{180}\psi^{(i)}(t)\right) - \\ &u^{(o)}(t)\cos\left(\frac{\pi}{180}\psi^{(o)}(t)\right). \end{aligned}$$

평가 단위는 자선-타선 쌍의 조우 구간이며, 각 조우 구간을 하나의 케이스로 정의하였다. 케이스 내에서 기준 시각 t_0 는 예측 지평선 $H = 6$ 분에서의 예측 최소 분리거리 $\hat{d}^{\min}(t)$ 를 기준으로 일관되게 선정하였다. 예측 최소 분리거리 $\hat{d}^{\min}(t)$ 는 일정 속도 가정 하에서 예측 지평선 H 내 최소 분리거리로 정의하며, 식 (8)와 같다. SOG는 kn 단위로 제공되며(1kn = 1NM/h), 이에 시간 변수 Δ 는 hour 단위로 두고 $H = 6\text{min} = 0.1h$ 로 사용하였다.

$$\begin{aligned} \hat{d}^{\min}(t) &= \\ \min_{0 \leq \Delta \leq H} &\sqrt{(\Delta x(t) + \Delta v_x(t)\Delta t)^2 + (\Delta y(t) + \Delta v_y(t)\Delta t)^2}. \end{aligned} \tag{8}$$

t_0 는 조우가 위험 수준으로 진입하는 시점을 일관된 기준으로 정하기 위해, $\hat{d}^{\min}(t)$ 의 임계 진입을 우선하는 규칙으로 결정하였다. 선택 규칙은 식 (9a)~(9b)와 같다.

$$\begin{aligned} T_{0.2} &= t \mid \hat{d}^{\min}(t) \leq 0.20, \\ T_{0.5} &= t \mid \hat{d}^{\min}(t) \leq 0.50, \end{aligned} \tag{9a}$$

$$t_0 = \begin{cases} \min T_{0.2}, & \mid T_{0.2} \mid > 0, \\ \min T_{0.5}, & \mid T_{0.2} \mid = 0, \mid T_{0.5} \mid > 0 \\ \hat{d}^{\min}(t), & \mid T_{0.5} \mid = 0 \end{cases} \tag{9b}$$

조우 구간 간 비교 조건을 동일하게 유지하기 위해 리플레이 윈도우는 식 (10)과 같이 고정하였다.

$$W = [t_0 - 5\text{min}, t_0 + 10\text{min}] \tag{10}$$

W 는 궤적 및 위험도 시계열을 동일 조건으로 기록·시각화하기 위한 공통 구간이며, 모델 출력 평가는 t_0 시각의 입력 s_t 에 대한 출력 y_{t_0} 를 기준으로 수행하였다.

3.2 비교 설정

제안 요소인 LoRA 어댑터와 RAG 모듈의 기여도를 분석하기 위해, 기본 모델(Base)을 기준으로 다음 네 가지 비교군을 구성하였다: B (Base), B+L (Base+LoRA), B+R (Base+RAG), B+L+R (Base+LoRA+RAG). 모든 비교군은 동일한 조우 입력 x_t 와 동일한 출력 포맷 y_t 를 공유하며, 오직 LoRA 및 RAG 적용 여부에서만 차이를 갖는다.

또한, 생성 기반 모델의 특성상 샘플링 확률 분포에 따라 결과가 달라질 수 있음을 고려하여, Nucleus Sampling의 Top-p 파라미터를 {0.2, 0.4, 0.6, 0.8, 1.0}으로 변화시키며 불확실성에 대한 민감도(Sensitivity)를 평가하였다. Top-p를 제외한 Temperature (0 고정) 등의 디코딩 설정은 모든 조건에서 동일하게 통제하였다.

3.3 평가 지표

본 연구는 제안 모델의 성능을 효율성, 안전성, 행동 성능의 세 가지 관점에서 종합적으로 평가하였다. 모든 지표는 개별 조

우 사례(case) 단위로 산출되었으며, 각 비교군 및 Top-p 설정에 따른 평균과 변동을 분석하였다. 구체적인 평가지표의 정의와 산출식은 다음과 같다.

첫째, final score(S_{final})는 모델의 종합적인 의사결정 효율성(Efficiency)을 나타내는 0~1 범위의 점수이다. 그래프에서 볼 수 있듯이 점수가 높을수록 우수한 성능을 의미하므로, 식 (11)과 같이 행동 정확도(S_{acc})와 내재적 안정성($1 - R_{illegal}$)의 가중 합으로 정의된다.

$$S_{final} = w_1 \cdot S_{acc} + w_2 \cdot 1 - R_{illegal} \quad (11)$$

여기서 w_1 과 w_2 각 지표의 중요도를 반영하는 가중치 계수이며, $R_{illegal}$ 은 후술할 내재적 위반율이다.

둘째, 내재적 위반율(Intrinsic Illegal Rate, $R_{illegal}$)은 안전 쉘드의 개입 없이 LLM이 생성한 원시 출력(Raw Output) 자체가 규정이나 물리적 안전 제약을 위반한 비율이다. 이는 전체 테스트 케이스 집합 N 에 대하여 식 (12)과 같이 계산된다.

$$R_{illegal} = \frac{1}{|N|} \sum_{i \in N} \mathbf{I}(y_{raw}^{(i)} \notin A_{valid}) \quad (12)$$

여기서 $y_{raw}^{(i)}$ 는 i 번째 사례에서의 LLM 원시 출력 행동, A_{valid} 는 규정상 허용되는 유효 행동 집합, $\mathbf{I}(\cdot)$ 는 지시 함수(Indicator function)이다. 그래프의 "Lower is Safer" 표시와 같이, 이 값이 0에 수렴할수록 모델이 내재적인 안전성을 갖추었음을 의미한다.

마지막으로 행동 F1 점수(action F1 score, $F1_{action}$)는 기준 레이블(reference label) 대비 행동 선택의 일치도를 나타내는 성능(performance) 지표로, 식 (13)과 같이 정밀도(precision)와 재현율(recall)의 조화 평균으로 산출된다.

$$F1_{action} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}} \quad (13)$$

4. 성능 분석

4.1 Top-p 스윙 기반 모델 비교

본 연구에서는 LLM의 생성 텍스트 다양성을 제어하기 위해 nucleus sampling (top-p) 기법을 적용하였다(Holtzman et al., 2020). Top-p 샘플링은 모델이 예측한 다음 토큰들의 확률 분포에서, 누적 확률(cumulative probability)이 p 값이 되는 상위 토큰들의 집합 $V^{(p)}$ 내에서만 샘플링을 수행하는 방식이다. 이는 다음 식 (14)와 같이 정의된다.

$$\sum_{x \in V^{(p)}} P(x|x_{1:i-1}) \geq p \quad (14)$$

여기서 p 는 생성의 무작위성을 결정하는 임계값이다. 자율운항 시나리오에서 p 값이 너무 낮으면 결정론적이고 반복적인 답변을 생성하여 복잡한 조우 상황에 유연하게 대처하지 못할 수 있으며, 반대로 너무 높으면 환각이나 일관성 없는 회피 행동을 유발할 위험이 있다. 따라서 본 연구에서는 p 값을 0.2에서 1.0까지 변화 시키며(sensitivity analysis), 모델의 추론 다양성이 항해 안전성과 효율성에 미치는 영향을 분석하였다.

구체적으로, 네 가지 비교 모델(B, B+L, B+R, B+L+R)을 대상으로 효율성(final score), 행동 성능(action F1 score), 안전성(intrinsic illegal rate) 지표를 측정하였으며, 그 결과는 Fig. 3-5에 제시하였다.

먼저 효율성 지표(Final Score)의 변화를 살펴보면(Fig. 3), Baseline(B)은 Top-p가 증가함에 따라 점수가 상승하는 경향을 보이나, 샘플링 설정 변화에 따른 변동성이 크게 나타났다. 반면 B+L+R은 중간 Top-p 구간에서 높은 점수 수준에 도달한 뒤, 고 Top-p 구간에서도 성능 저하 없이 비교적 안정적으로 유지되는 양상을 보였다. 이는 RAG가 판단을 규정 텍스트에 정렬시켜 생성의 탐색 공간을 제한하고, LoRA가 구조화 출력 및 규칙 선택 패턴을 일관되게 따르도록 보조함으로써, 샘플링 불확실성 변화가 효율 지표의 급격한 흔들림으로 전이되는 것을 효과적으로 완화한 결과로 해석된다.

행동 성능 지표(action F1 score)의 결과는 Fig. 4와 같다. 전반적으로 top-p가 증가할수록 샘플링 다양성이 확대되어 상황별 행동 선택이 개선되며 F1 점수가 상승하는 경향이 관찰된다. baseline(B)은 top-p \approx 1.0 구간에서 B+L+R과 유사한 수준에 도달하지만, 낮은 top-p 구간에서는 출력이 제한되면서 행동 선택이 보수적(또는 passive)으로 치우쳐 F1 저하가 나타난다. 반면 B+L+R은 낮은 top-p에서도 규칙 적용과 행동 선택이 비교적 안정적으로 유지되어 설정 변화에 대한 강건성이 높다. 또한 Fig. 3의 효율 점수 곡선이 Fig. 4와 유사한 추세를 보이는데, 적절한 행동 선택(F1 향상)이 불필요한 회피·감속을 줄여 운항 효율 역시 함께 개선되기 때문이다.

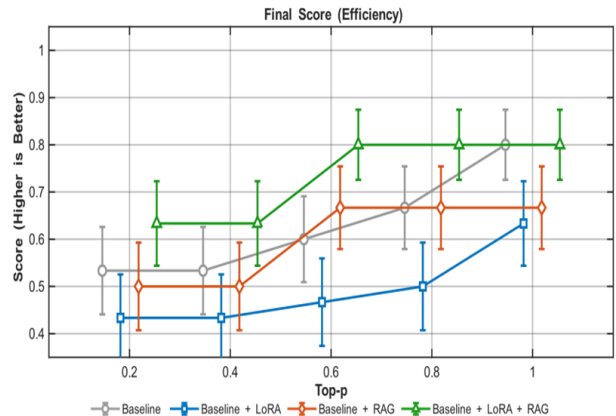


Fig. 3 Trends in final efficiency scores with respect to Top-p sampling hyperparameters

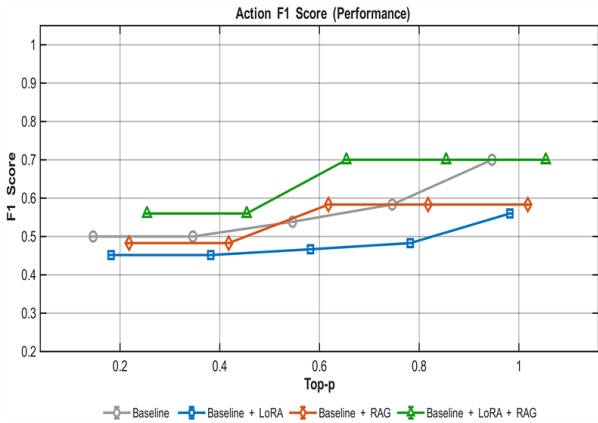


Fig. 4 Analysis of action performance improvement with increasing sampling diversity

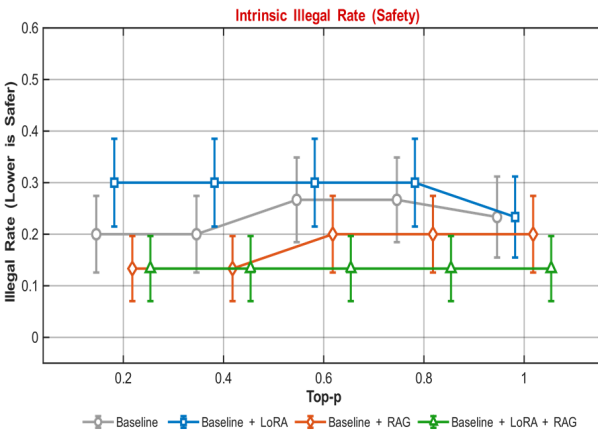


Fig. 5 Comparison of intrinsic illegal rates across model configurations and sampling ranges

안전성 지표(intrinsic illegal rate)는 Fig. 5에 제시하였다. B와 B+R은 top-p 변화에 따라 성능과 안전 간 상충 관계(Trade-off)가 관찰되며, B+L은 전 구간에서 상대적으로 높은 위반율을 보였다. 반면 제안 모델(B+L+R)은 top-p 전 구간에서 낮은 위반율을 일관되게 유지하였다. 이는 RAG가 규정 오인의 가능성을 줄이고, LoRA가 규칙 선택·조치 제안의 형식적 일관성을 강화함으로써, 샘플링 불확실성 증가 상황에서도 위반 행동으로의 이탈을 억제하는 방향으로 작동했음을 시사한다. 특히 Fig. 5의 낮은 내재적 위반율은 결과적으로 안전 쉼드의 개입 빈도($Z_t = 1$)가 낮음을 의미한다. 이는 제안 모델이 사후 보정(Safety Shield)에 전적으로 의존하는 것이 아니라, LLM 추론 단계에서부터 규정에 부합하는 안전한 경로를 생성하고 있음을 시사한다. 이어지는 4.2 절에서는 이러한 판단이 실제 시나리오에서 어떻게 구체화되는지 확인한다.

이러한 분석 결과를 종합하여 가장 높은 다양성 설정(top-p=1.0)에서의 모델별 정량 성능을 요약하면 Table 3과 같다. 제안 모델(B+L+R)은 종합 점수(Final Score) 0.80, 행동 정확도(Action F1) 0.71을 기록하여 모든 지표에서 가장 우수한 성능을 보였다. 특히 안전성(safety) 지표에서는 0.13의 가장 낮은 위반율을 기록

Table 3 Summary of quantitative performance metrics across comparative models

Model	Final Score	Action F1	Safety
Base	0.65	0.68	0.23
Base + LoRA	0.62	0.55	0.30
Base + RAG	0.67	0.58	0.20
Base + LoRA + RAG	0.80	0.71	0.13

하여, Baseline(0.23) 대비 약 43%의 개선 효과를 확인하였다.

4.2 AIS 기반 정성 사례 분석

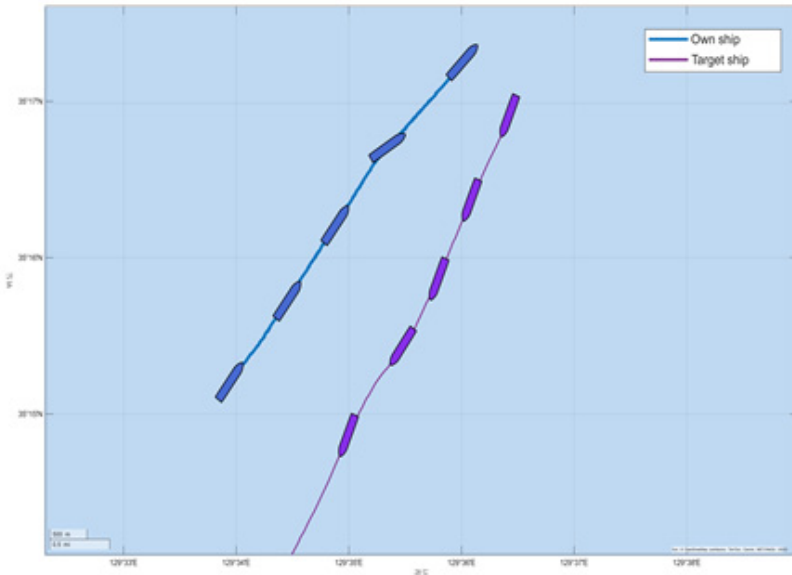
본 절에서는 앞선 4.1절의 정량 평가가 모델의 통계적 성능을 입증했다면, 본 절에서는 실제 해상 시나리오에서 시스템이 어떠한 논리적 인과관계를 거쳐 의사결정에 도달하는지 심층적으로 분석한다. Fig. 6~9는 실제 AIS 항적 데이터를 기반으로, 기하학적 위험 인지에서 시작하여 RAG 기반 규정 해석과 LLM 행동 생성을 거쳐 최종적인 안전 쉼드 검증에 이르는 일련의 처리 과정을 단계별로 도시하였다. 이를 통해 수치만으로는 파악하기 어려운 '설명 가능한 판단 근거(Explainability)'와 '결정론적 안전 보장(Safety Assurance)'의 상호작용 메커니즘을 구체적으로 확인한다.

Fig. 6은 정면 조우(Head-on) 상황에 대한 분석 결과를 보여준다. 시스템은 양 선박이 정면으로 마주 보는 형세와 좁혀지는 거리를 근거로 충돌 위험을 식별하고, 해당 조건이 COLREGs 제 14조 적용 대상임을 바탕으로 우현 변칙을 조치로 제시하였다. 이후 안전 쉼드는 제안된 조치가 물리적 안전 기준(DCPA/TCPA)을 만족하는지 검증한 뒤 최종 행동으로 확정하였다.

Fig. 7은 교차 조우(Crossing) 상황에서 피항 의무를 이행한 사례를 제시한다. 자선(파란색) 기준으로 타선(보라색)이 우현 방향에서 횡단(Starboard-to-Port crossing)하며 접근함에 따라 자선이 피항선(give-way vessel)으로 판단된다. 시스템은 COLREGs 제15조-제16조에 근거해 선수 횡단을 피하고 선미 통과(pass astern)를 유도하는 침로 변경(heading alteration) 조치를 도출하였으며, 우측 패널의 최종 명령과 같이 실행 단계로 확정하였다. 이후 안전 쉼드 검증을 통해 해당 조치가 물리적 안전 기준(DCPA/TCPA)을 만족함을 확인한 뒤 최종 행동으로 적용하였다.

Fig. 8은 속도 제어를 포함한 복합적 판단(선미 통과 전략)의 사례를 보여준다. 단순 침로 유지 시 상대선의 진로를 가로지를 위험이 존재하므로, 시스템은 무리한 횡단 대신 COLREGs 제8조에 근거하여 필요 시 감속과 침로 변경을 병행하는 전략을 선택하였다. 특히 Fig. 8의 우측 패널은 LLM 입력 프로세스의 구체적인 예시로, 센서로부터 수집된 기하 정보(panel 1)와 RAG를 통해 검색된 COLREGs 제8조/제15조(panel 4)가 작업자의 수동 개입 없이 시스템에 의해 자동으로 직렬화(serialization)되어 프롬프트로 변환되는 과정을 나타낸다.

이러한 전략이 시간축에서 수행되는 과정은 Fig. 10에 도시되



1. SENSOR DATA & GEOMETRY

SITUATION: Head-on (Rule 14) - **CONFIRMED**

- TARGET: Reciprocal Course (Delta COG > 160 deg)
- ASPECT: Mutual Dead-ahead (RBoW < 15 deg AND RBoTarget < 15 deg)
- RANGE: 0.12 NM (222 m) - **[CLOSE-QUARTERS]**
- DCPA: 0.11 NM (204 m) - **[CRITICAL]**

2. OOW APPRAISAL

OBSERVATION: Constant bearing with decreasing range.

ASSESSMENT: Risk of collision cannot be ruled out.
(Condition: $d(Brg)/dt < Epsbeta$ and $TCPA < Tmin$)

RATIONALE: In a head-on situation, both vessels are required to alter course to starboard to pass port-to-port. (See Citation [E1], [E2]).

3. HYBRID DECISION & EXECUTION

[STEP 1] LLM PROPOSAL
> Action: Alter course to Starboard (Substantial).

[STEP 2] SAFETY SHIELD VERIFICATION:
> Check: (Pred. DCPA \geq $dmin$) AND (Pred. TCPA \geq $tmin$)
> Result: **PASSED** (Intervention = 0)

[STEP 3] FINAL COMMAND (Rules 14, 8, 34):
> HELM: HARD STARBOARD (Delta Psi \geq 20 deg).
> ENGINE: Standby
> SIGNAL: One short blast (intention).

4. RAG EVIDENCE (Retrieved Context)

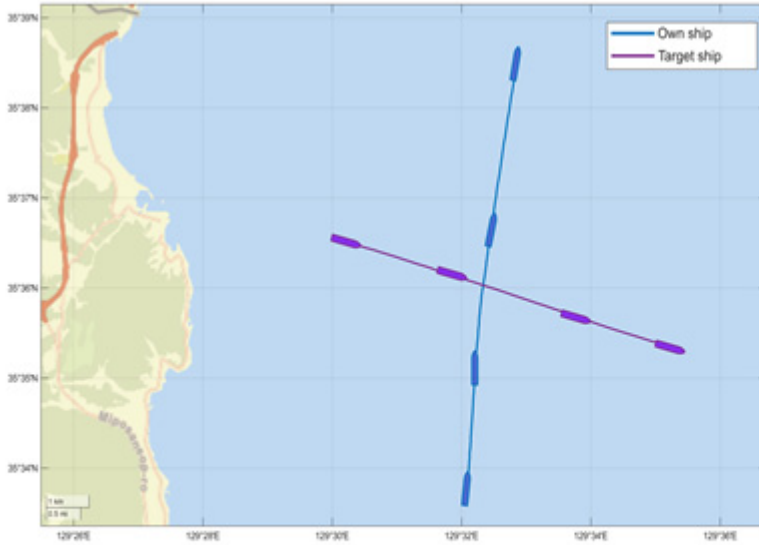
[E1] COLREGs Rule 14 (Head-on): "...each shall alter course to stbd..."
 [E2] COLREGs Rule 8 (Action): "...action shall be positive, early..."
 [E3] COLREGs Rule 34 (Sound): "...one short blast to mean I am altering..."
 (Source: IMO Convention 1972, Consolidated 2018)

Fig. 6 Qualitative result for Head-on encounter

었다. 그래프는 위에서부터 최종 행동 인덱스와 안전 쉴드 활성화 ($\sigma_t = 1$), 제어 입력(SOG/COG), 그리고 위험 지표(DCPA/TCPA 및 위험 구간)를 시계열로 도시하였다. Fig. 9(a)에서 위험 구간 ($f_{risk} = 1$)진입 예측 시 안전 쉴드가 활성화($\sigma_t = 1$)되어 개입하며, 이후 감속 및 침로 변경이 적용됨에 따라 DCPA/TCPA가 임계 조건을 만족하는 방향으로 전이되어 위험 상태가 해소됨을 확인할 수 있다. 이를 통해 정성적으로 규정된 “충분히 이른 회피” 요구가 데이터 기반의 자동화된 프롬프트를 통해 구체적 제어 명령으로 변환되고, 최종적으로 안전 쉴드(σ_t)의 물리적 검증을 거

쳐 행동이 확정되는 과정을 입증한다.

Fig. 9은 복수 선박이 얽힌 고난도 조우 상황에서의 충돌 위험 해소 과정을 나타낸다. 정면 및 교차 위험이 혼재된 상황에서 단순한 좌현 변침은 특정 선박(T2)과의 충돌을 유발할 수 있으므로, 시스템은 상충 관계를 분석하여 T1-T3 모두에 대해 동시에 안전 영역을 확보할 수 있는 우현 전타(Hard Starboard)를 최종 행동으로 도출하였다. 이는 다선박 상황에서 국소적으로 합리적인 조치가 전체 안전을 훼손할 수 있음을 고려하여, 단일 파이프라인 내에서 규칙 적용과 물리적 검증을 결합해 일관된 최종 행동을 산출했음을 보여준다.



1. SENSOR DATA & GEOMETRY

SITUATION: Crossing (Rule 15) - **CONFIRMED**

- TARGET STATUS: Crossing from Starboard to Port
- RELATIVE BEARING: Starboard Bow (+045 deg)
(Target is in the Sector: 005 ~ 112.5 deg)
- RANGE: 0.30 NM (556 m) - **[RISK OF COLLISION]**
- DCPA: 0.15 NM (278 m) - **[CRITICAL]**

2. OOW APPRAISAL

OBSERVATION: Target is on own starboard side.
Bearing is constant, range is closing.

ASSESSMENT: Own ship is the **GIVE-WAY** vessel.
(Condition: Target in 0~112.5 deg RB AND Risk exists)

RATIONALE: Per Rule 15, when two power-driven vessels are crossing, the vessel which has the other on her own starboard side shall keep out of the way and avoid crossing ahead.
(See Citation [E1], [E2]).

3. HYBRID DECISION & EXECUTION

[STEP 1] LLM PROPOSAL
> Action: Alter course to Starboard to pass astern.

[STEP 2] SAFETY SHIELD VERIFICATION:
> Check: (Pred. DCPA >= dmin) AND (No Crossing Ahead)?
> Result: **PASSED** (Intervention = 0)

[STEP 3] FINAL COMMAND (Rules 15, 16, 8):
> HELM: STARBOARD 20 (Aiming to pass target stern).
> ENGINE: Maintain or Slow down if necessary.
> SIGNAL: One short blast (intention).

4. RAG EVIDENCE (Retrieved Context)

[E1] COLREGs Rule 15 (Crossing): "...vessel which has the other on her own starboard side shall keep out of the way..."
 [E2] COLREGs Rule 16 (Action by Give-way): "...take early and substantial action to keep well clear."
 [E3] COLREGs Rule 8 (Action): "...avoid crossing ahead..."
 (Source: IMO Convention 1972, Consolidated 2018)

Fig. 7 Qualitative result for crossing encounter

또한 Fig. 11은 동일 사례에 대해 최종 행동 인덱스와 쉴드 활성화(σ_t), 제어 입력(SOG/COG), 그리고 DCPA/TCPA 기반 위험 지표의 시간 변화를 함께 제시하여, hard starboard가 선택·확정되고 위험이 완화되는 과정을 시간축에서 확인할 수 있도록 보장한다. 이상의 사례는 기하학적 위험 인지가 규정 검색 및 규칙 적용으로 연결되고, 이것이 다시 안전 쉴드(σ_t)의 물리적 검증을 거쳐 최종 행동으로 확정되는 일련의 과정을 사례 수준에서 확인한다.

이상의 사례는 기하학적 위험 인지가 규정 검색 및 규칙 적용으로 연결되고, 이것이 다시 안전 쉴드의 물리적 검증을 거쳐 최

종 행동으로 확정되는 일련의 과정을 사례 수준에서 확인한다.

5. 결론

본 논문은 실제 AIS 데이터 기반 조우 상황(오프라인 리플레이)에서 COLREGs 준수 의사결정을 설명 가능한 형태로 생성하기 위해, RAG 기반 근거 검색과 LLM 기반 판단 생성을 결합한 하이브리드 의사결정 프레임워크를 제안하였다. 본 연구는 선박을 직접 제어하는 경로·제어 알고리즘을 제시하기보다, 조우 상황을 해석하고 적용 규칙을 선택하며 그 근거를 인용해 설명하는



1. SENSOR DATA & GEOMETRY

SITUATION: Crossing (Rule 15) - **CONFIRMED**

- TARGET STATUS: Crossing (Target on own starboard side)
- ASPECT: Target on Starboard Bow (+045 deg)
(Target is in the Sector: 15 ~ 112.5 deg RB)
- RANGE: 0.35 NM (648 m) - **[CLOSE-QUARTERS]**
- DCPA: 0.15 NM (278 m)

2. OOW APPRAISAL

OBSERVATION: Target is on own starboard side.
Bearing is constant, range is closing.

ASSESSMENT: Own ship is the **GIVE-WAY** vessel.
Condition: $(15 \text{ deg} \leq \text{RB}_{\text{own}} < 112.5 \text{ deg})$ AND
 $(\text{DCPA} < d_{\text{min}})$ AND $(\text{TCPA} < t_{\text{min}})$

RATIONALE: Per Rule 15, vessel with the other on her starboard side shall keep out of the way and avoid crossing ahead.
(See Citation [E1], [E2]).

3. HYBRID DECISION & EXECUTION

[STEP 1] LLM PROPOSAL (a):
> Action: Alter course to Starboard to pass astern.

[STEP 2] SAFETY SHIELD VERIFICATION
(Pred = Prediction after applying LLM Action)
> Check: $(\text{Pred. DCPA} \geq d_{\text{min}})$ AND
 $(\text{Pred. TCPA} \geq t_{\text{min}})$ AND (Passing astern)?
> Result: **PASSED** (Intervention = 0)

[STEP 3] FINAL COMMAND (Rules 15, 16, 8, 34):
> HELM: STARBOARD 20 (Aiming to pass target stem).
> ENGINE: Maintain or Slow down if necessary.
> SIGNAL: One short blast (intention).

4. RAG EVIDENCE (Retrieved Context)

[E1] COLREGs Rule 15 (Crossing): "...vessel which has the other on her own starboard side shall keep out of the way..."
[E2] COLREGs Rule 16 (Action by Give-way): "...take early and substantial action to keep well clear."
[E3] COLREGs Rule 8 (General Action): "Action shall be positive, made in ample time, and with due regard to good seamanship."
[E4] COLREGs Rule 34 (Sound): "...one short blast..."
(Source: IMO Convention 1972, Consolidated 2018)

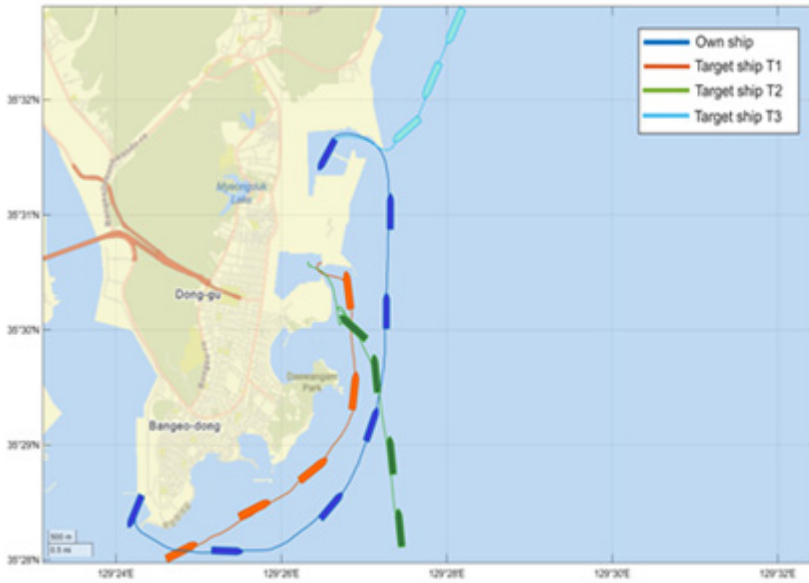
Fig. 8 Strategic Decision-Making for pass astern

“항해사 관점의 판단 레이어” 구축에 초점을 두었다. 이를 위해 COLREGs 원문과 공신력 있는 해설을 검색·참조하여 판단이 규정 텍스트에 기반하도록 유도하였고, LLM이 규칙 선택과 조치 제안을 근거 인용과 함께 구조화된 출력으로 생성하도록 설계하였다.

또한, LLM이 물리적 위험을 과소평가하거나 기하학적 추론 오류를 범할 수 있는 실패 모드를 고려하여, CPA/TCPA 기반의 결정론적 안전 쉼드를 런타임 안전 계층으로 도입하였다. 이에 따라 규칙의 해석과 판단은 LLM이 수행하되, 최종 행동의 물리

적 안전은 결정론적 모듈이 강제하는 역할 분리(Decoupling) 구조를 통해 안전 필수 운항 시스템으로의 적용 가능성을 제시하였다.

정량 평가에서는 top-p 스위치를 통해 샘플링 설정 변화에 대한 성능 및 안전 민감도를 비교하였다. 그 결과 제안 방법은 top-p가 증가하는 조건에서도 안전 지표의 악화를 상대적으로 억제하는 경향을 보였다. 이는 RAG가 생성 과정에서 참조 가능한 근거를 제공하여 규정 적용의 임의성을 줄이고, 안전 쉼드가 기하학적 안전 제약을 만족하지 않는 조치를 런타임에서 차단·대체함으로



1. MULTI-TARGET AWARENESS

SITUATION: 3-Vessel Encounter (High Complexity)

- [T1] Head-on (Red): Range 0.30 NM -> **[MDT: PRIORITY 1]**
(Immediate Danger, Dead Ahead)
- [T2] Crossing (Yellow): Range 0.80 NM -> **[PRIORITY 2]**
(Approaching from Stbd Bow)
- [T3] Head-on (Sky Blue): Range 1.20 NM -> **[PRIORITY 3]**
(Approaching from North Dead Ahead)

3. HYBRID DECISION & EXECUTION

- [STEP 1] LLM PROPOSAL (Unified Action):
> Action: Alter course to Starboard (Substantial, >30 deg).
- [STEP 2] SAFETY SHIELD VERIFICATION:
> Check: Risk(T1) & Risk(T2) & Risk(T3) < Safe
> Result: **PASSED** (Intervention = 0)
- [STEP 3] FINAL COMMAND (Rules 14, 15, 8, 2):
> HELM: HARD STARBOARD (To clear entire sector).
> ENGINE: Standby / Slow down if turning radius is tight.
> SIGNAL: One short blast (intention).

2. OOW APPRAISAL (Integrated Assessment)

- OBSERVATION: Dense traffic situation.
- T1 & T3 (Dead-Ahead): Both require Starboard Turn (Rule 14).
- T2 (Stbd Side): Requires Give-way (Starboard Turn, Rule 15).
- CONFLICT CHECK: Single maneuver for all targets?
> Analysis: Starboard turn satisfies Rules for T1, T2, and T3.
> (Left turn is **FATAL** due to T2).
- RATIONALE: A large starboard alteration is the only viable option to clear T1/T3 (Port-to-port) and pass astern of T2.

4. RAG EVIDENCE (Retrieved Context)

- [E1] Rule 14 (Head-on): "Each shall alter course to stbd..."
- [E2] Rule 15 (Crossing): "Keep out of the way... avoid crossing ahead."
- [E3] Rule 8 (Action): "Action shall be positive... ample time..."
- [E4] Rule 2 (Responsibility): "Special circumstances..."
(Source: IMO Convention 1972, Consolidated 2018)

Fig. 9 Resolution of Multi-Vessel conflict

써, 생성 불확실성이 안전 저하로 직접 전이되는 것을 완화하기
때문으로 해석된다. 즉, 본 프레임워크는 확률적 생성 모델의 불
확실성 하에서도 성능과 안전 간 상충 관계(trade-off)를 완화할
수 있는 구조적 장치를 제공한다.

정성 평가에서는 AIS 리플레이 기반 대표 조우 사례를 통해 기
하 요약, 규칙 선택, 근거 인용, 그리고 안전 실드 검증이 하나의
파이프라인으로 연결됨을 제시하였다. 다선박, 정면, 교차 조우
등 복합 상황에서 시스템이 도출한 판단의 근거와 안전 검증 과
정이 함께 제시되어, 항해사 관점의 "상황 인지-평가-조치" 흐름
이 실제 데이터 기반 시나리오에서도 구현 가능함을 확인하였다.

본 연구는 실제 AIS 데이터에서 충돌 임박 사례가 희소하여 위
험 구간 위주의 오프라인 검증을 수행했다는 한계를 갖는다. 향
후 연구에서는 다선박 상황의 확장, 다양한 해역 및 공개 AIS 데
이터에 대한 추가 검증, 근거 인용의 정확도에 대한 정량 평가,
그리고 경로-제어 계층과의 통합을 포함한 실시간 운용 환경에서
의 종단 간(end-to-end) 검증을 수행할 계획이다.

후 기

이 논문은 국립부경대학교 자율창의기술연구비(2025년)에 의

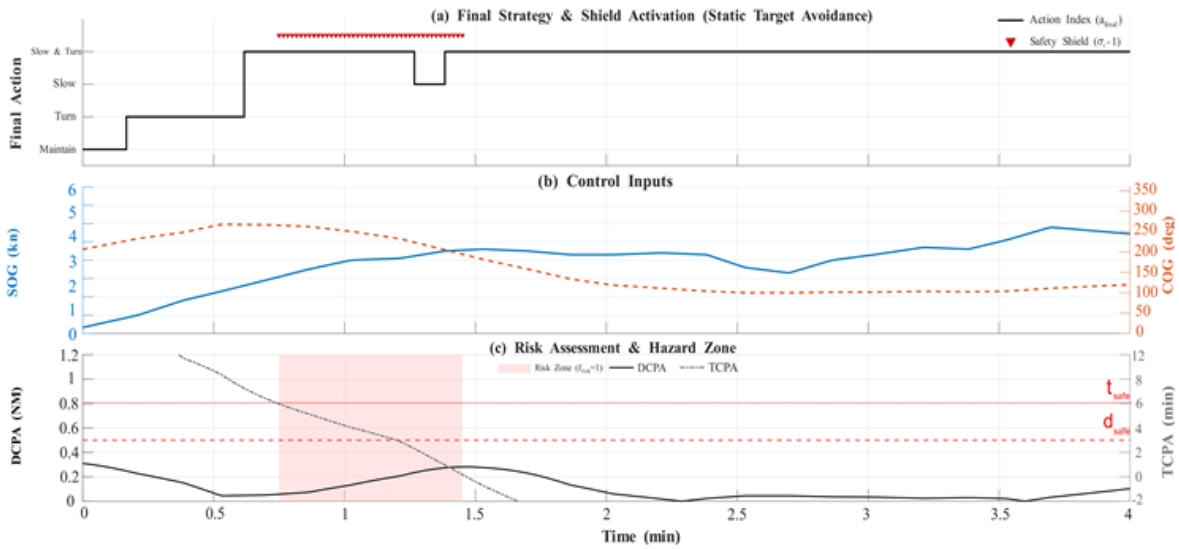


Fig. 10 Time-Series Visualization of Pass-Astern Strategy with Safety-Shield Activation

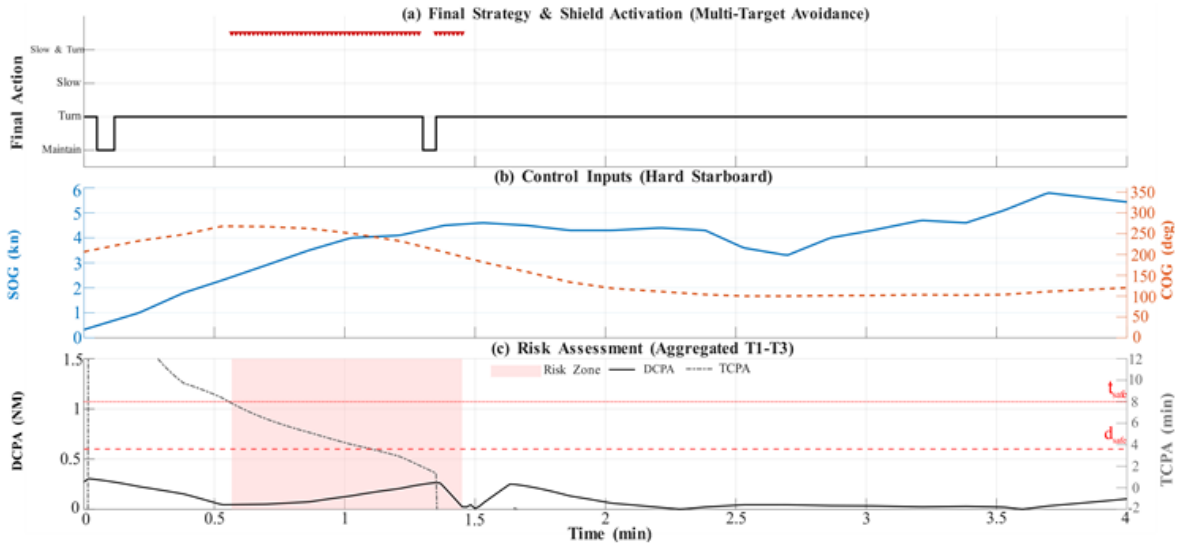


Fig. 11 Time-Series Analysis of Multi-Vessel Conflict Resolution with Safety-Shield Activation

하여 연구되었음. 또한, 2025년도 해양수산부 재원으로 해양수산 과학기술진흥원의 지원을 받아 수행된 연구임(RS-2024-00410 200, 한-유럽 첨단 해양모빌리티 연구거점 구축 및 공동연구).

References

Allianz Global Corporate & Specialty, 2025. Safety and shipping review 2025. Allianz Global Corporate & Specialty.
 Fiorini, P. and Shiller, Z., 1998. Motion planning in dynamic environments using velocity obstacles. *International Journal of Robotics Research*, 17(7), pp.760-772.
 Holtzman, A., Buys, J., Du, L., Forbes, M. and Choi, Y., 2020. The curious case of neural text degeneration.

International Conference on Learning Representations (ICLR), Addis Ababa, Ethiopia.
 Hu, E.J., Shen, Y., Wallis, P., Allen-Zhu, Z., Li, Y., Wang, S., Wang, L. and Chen, W., 2021. LoRA: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*.
 International Maritime Organization (IMO), 1972. Convention on the international regulations for preventing collisions at sea, 1972 (COLREGs). IMO: London.
 Ji, Z., Lee, N., Frieske, R., Yu, T., Su, D., Xu, Y., Ishii, E., Bang, Y.J., Madotto, A. and Fung, P., 2023. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12), pp.1-38.
 Johansen, T.A., Perez, T. and Cristofaro, A., 2016. Ship

collision avoidance and COLREGS compliance using simulation-based control behavior selection with predictive hazard assessment. *IEEE Transactions on Intelligent Transportation Systems*, 17(12), pp.3407-3422.

Lewis, P., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W. and Rocktäschel, T., 2020. Retrieval-augmented generation for knowledge-intensive NLP tasks. *Advances in Neural Information Processing Systems*, 33, pp.9459-9474.

Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., Rodriguez, A., Joulin, A., Grave, E. and Lample, G., 2023. LLaMA: open and efficient foundation language models. *arXiv preprint arXiv:2302.13971*.

Vagale, A., Bye, R.T., Oucheikh, R., Osen, O.L. and Fossen, T.I., 2021. Path planning and collision avoidance for autonomous surface vehicles II: A comparative study of algorithms. *Journal of Marine Science and Technology*, 26, pp.1307-1323.

Zhang, S., Chen, D., Xiao, Z., Lü, R., Liu, M., Wu, M. and Wang, Y., 2023. Language MPC: Large language models as decision makers for autonomous driving. *arXiv preprint arXiv:2310.03026*.

Authorship Contribution Statement

Jin-Hyeok Seo: Conceptualization, Data curation, Formal analysis, Investigation, Methodology, Software, Validation, Visualization, Writing – original draft, Writing – review & editing; **Seung-hyeon Lim:** Visualization, Investigation; **Seong-Hyeon Jeong:** Data curation, Formal analysis, Validation; **Sun-Hyuck Im:** Data curation, Formal analysis, Validation, Methodology; **Si-Won Kim:** Data curation, Methodology; **Yeon-Soo Kim:** Data curation, Visualization, ; **Jeong-Hyeon Kim:** Visualization; **Jong-Yong Park:** Conceptualization, Project administration, Supervision, Writing – review & editing.

