

빠르고 똑똑하면서 “신뢰가능한” 안전해석 AI 조교

2026. 05. 06

2026 한국원자력학회 춘계학술발표회
원자력 인공지능 강습회

Kyung Mo Kim, Ph.D.

Korea Institute of Energy Technology (KENTECH)
Department of Energy Engineering

kmokim@kentech.ac.kr

AI in Everywhere

Photonic chips



XR



Energy



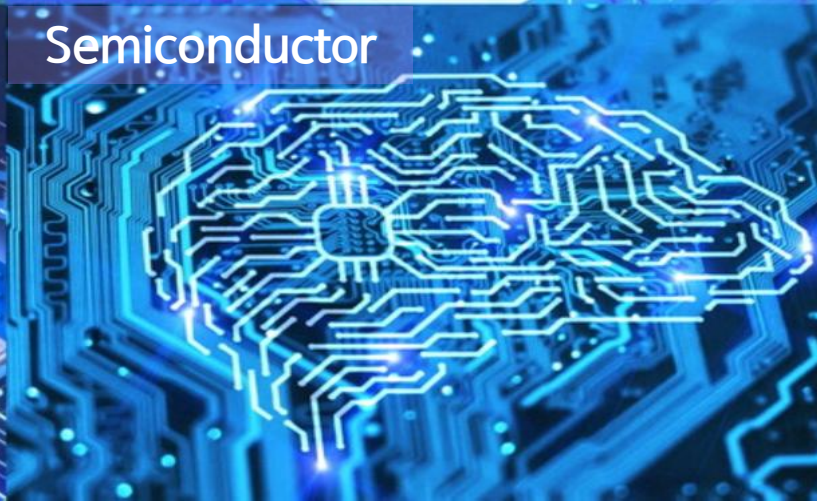
Perceptive soft robotics



Integrated computing



Semiconductor



High-precision medicine



인공지능 기술 발전

□ Artificial general intelligence (AGI) – 초지능 개발

nature

View all journals Search Log in

Explore content About the journal Publish with us Subscribe Sign up for alerts RSS feed

nature > news.feature > article

NEWS FEATURE | 30 December 2025

Science in 2050: the future breakthroughs that will shape our world – and beyond

Nuclear fusion. People on Mars. Artificial general intelligence. These are just some of the advances that could come by the mid-century mark.

By David Adam



Generative AI

AI Agents

Agentic AI

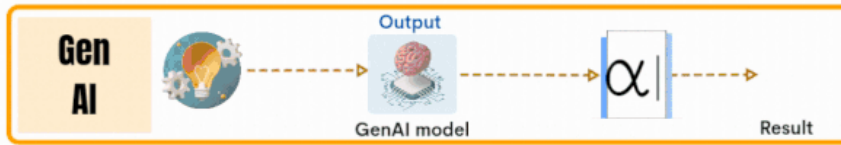


α | be the alpha
Your Intelligence Engine

ASPECT	GENERATIVE AI	AI AGENTS	AGENTIC AI
Purpose	Produces content like text or images.	Automates tasks using rules or behaviors.	Acts autonomously, making complex decisions.
Functionality	Focuses on creative and novel outputs.	Executes predefined tasks and actions.	Interacts and adapts dynamically to environment.
Examples	GPT-3, DALL-E	Chatbots, virtual assistants	Autonomous vehicles, intelligent robots
System Interaction	Generates output, not typically system-centric.	Functions within system guidelines.	Deeply embedded, continuously interacting with systems.
Learning Capability	Updated through retraining to enhance creativity.	Limited learning potential unless integrated with algorithms.	Designed to learn and adapt through interaction.
Adaptability	Improves creative output with updates.	Adaptable only within predefined constraints.	Highly adaptable, improving through experience.



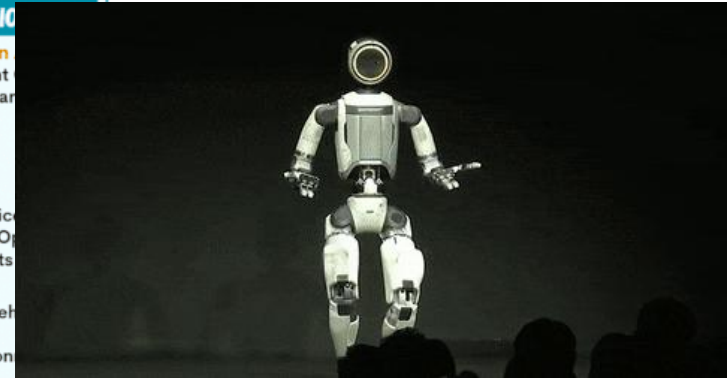
Hyundai Motor Group | Boston Dynamics



USE CASES APPLICATIONS

- Generative AI (Gen AI)**
 - Art and Content
 - Virtual Worlds and Gaming
 - Marketing and Advertising
- AI Agents**
 - Customer Service
 - Automation in Operations
 - Smart Assistants

- Agentic AI**
 - Autonomous Vehicles
 - Robotics
 - Dynamic Environment Management



인공지능에게 기대하는 바

□ 기존 계산/해석/예측의 한계를 극복할 수 있는 “정확성”

상관식/모델 한계

Bowring (1972)

$$q''_c = \frac{A + B(h_f - h_{in})}{C + L}$$

Katto & Ohno (1984)

$$q''_c = q''_{co} \left(1 + K \frac{h_f - h_{in}}{h_{fg}} \right)$$

PG-CHF (1986)

$$R = \frac{k_1 F_g}{f(P_r)(dT_r)^{k_2}} \frac{f(P_r, G)f(P_r, X_i)}{f_a} = f(P_r, G, X_i, X, q'', dT_r, F_g)$$

Shah (1987)

$$Bo = Bo_{UCC} \text{ for } Y \leq 10^6$$

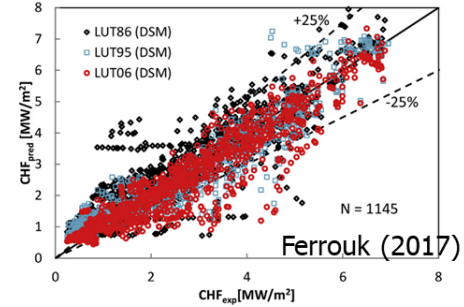
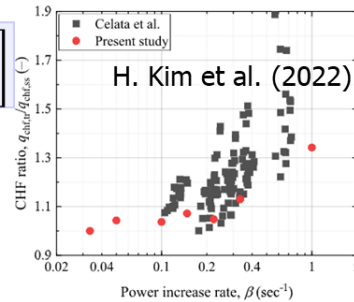
$$Bo = \min\{Bo_{UCC}, Bo_{LCC}\} \text{ for } Y > 10^6$$

$$q''_c = Bo \cdot G \cdot h_{fg}$$

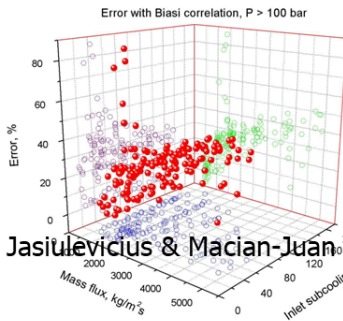
Biasi (1967)

$$q''_{CHF} = \frac{1.883 \times 10^7}{G^{1/6} D_h^n} \left[\frac{f_p}{G^{1/6}} - X_e \right]$$

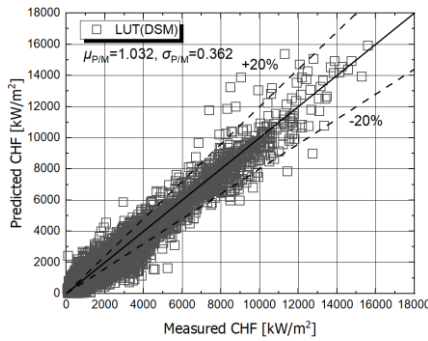
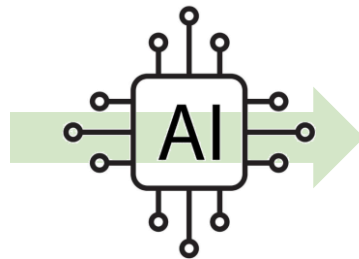
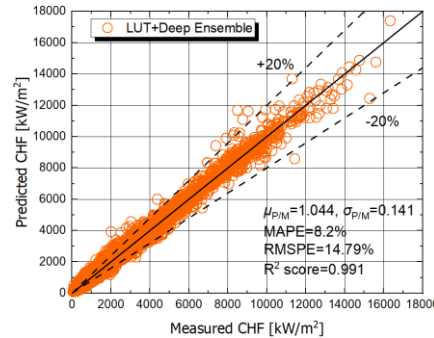
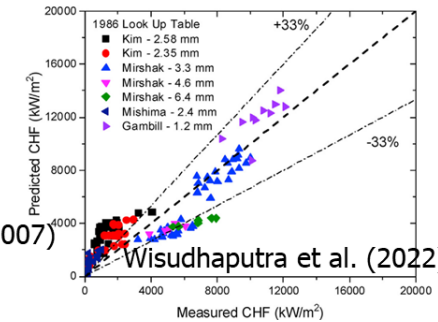
$$q''_{CHF} = \frac{3.78 \times 10^7}{G^{0.6} D_h^n} h_p [1 - X_e]$$



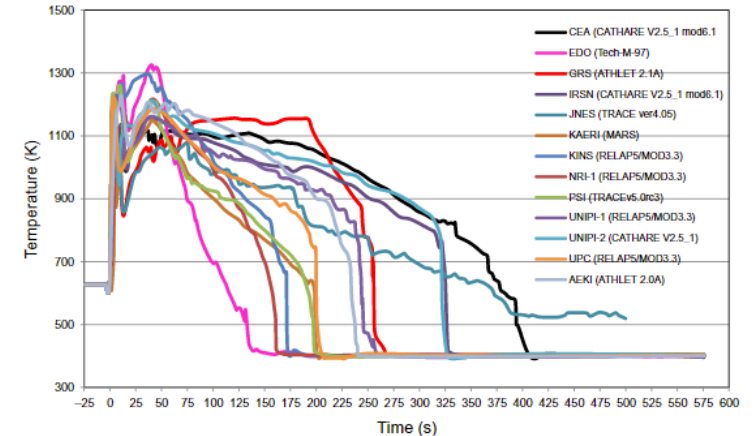
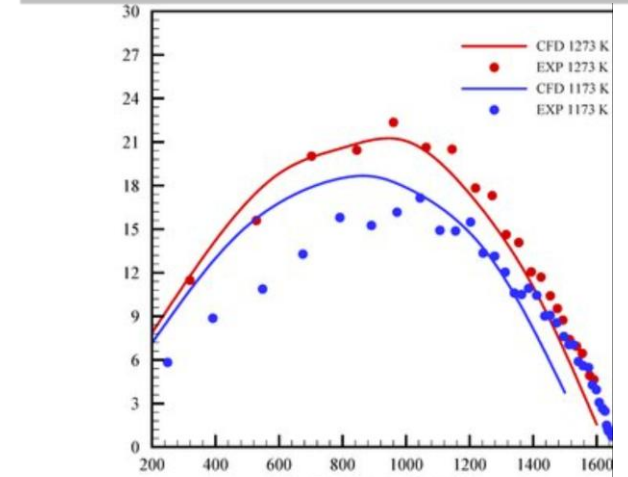
Geometry, i.e., Bundle, rectangular, etc.



Jasilevicius & Macian-Juan (2007)

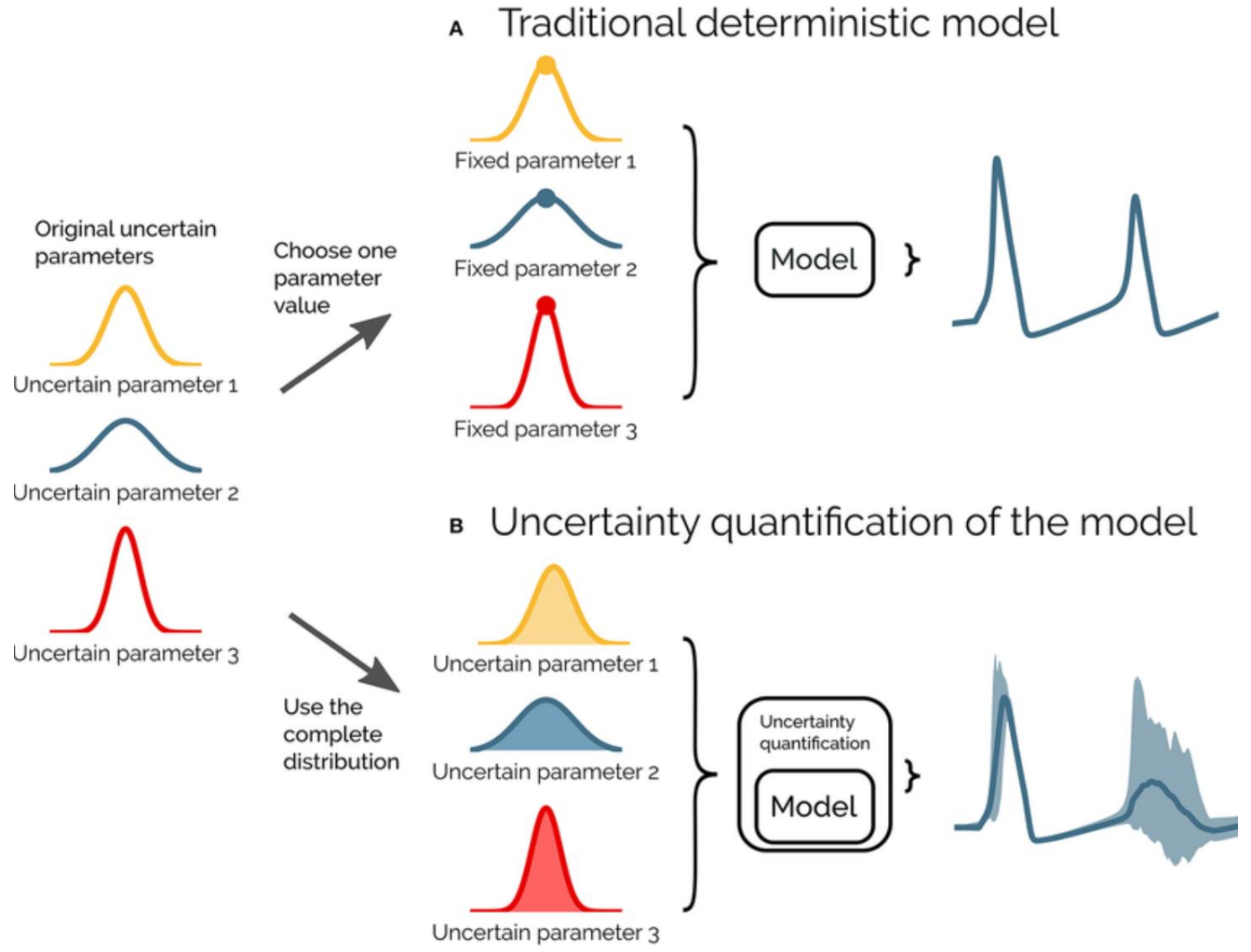
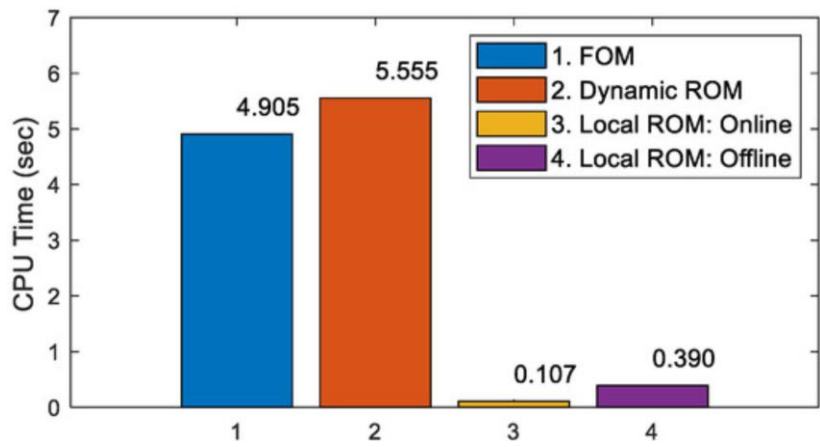
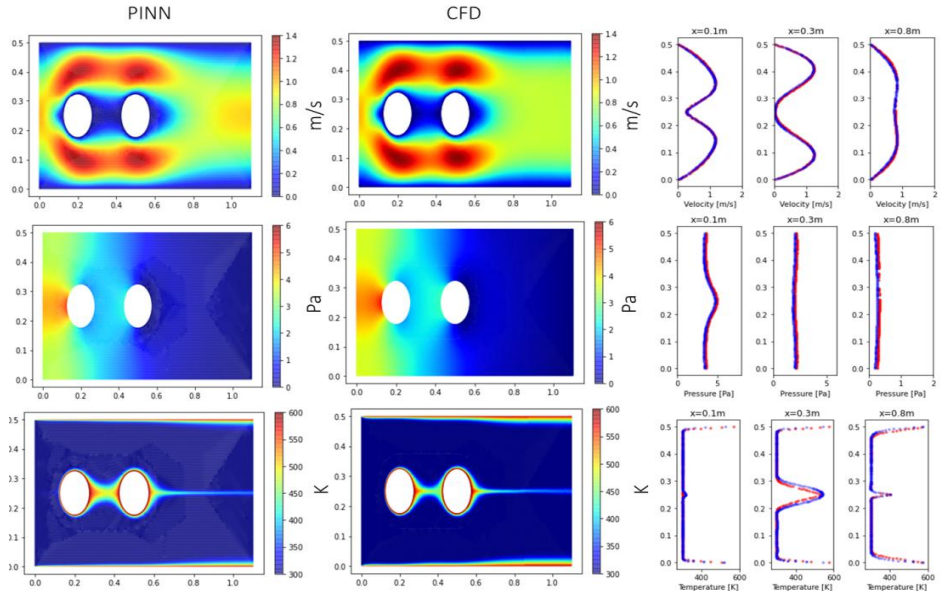


수치해석 한계



인공지능에게 기대하는 바

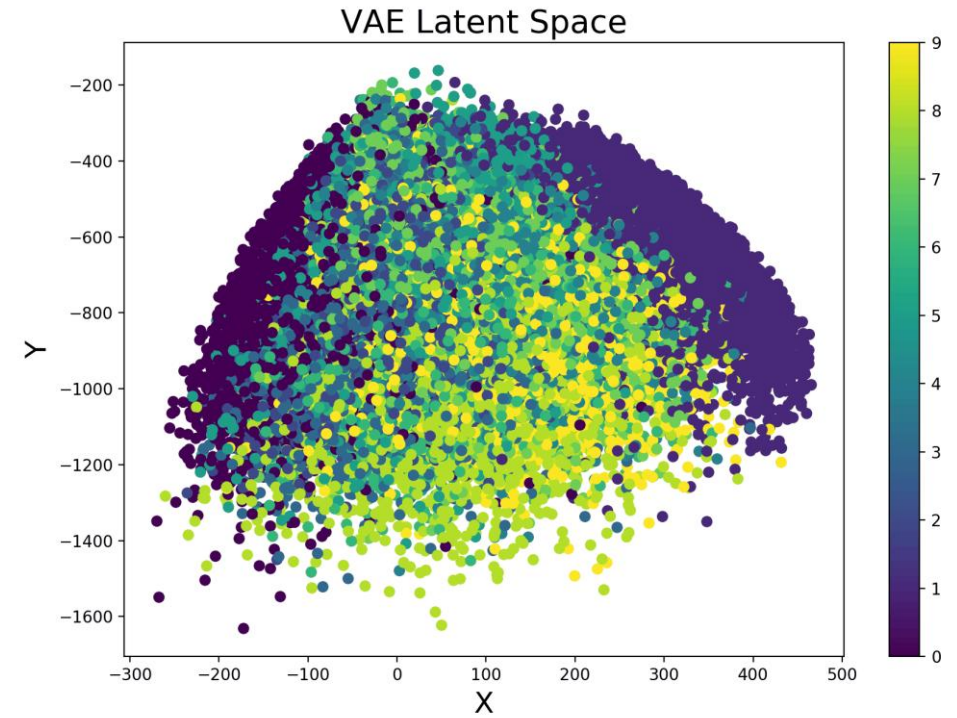
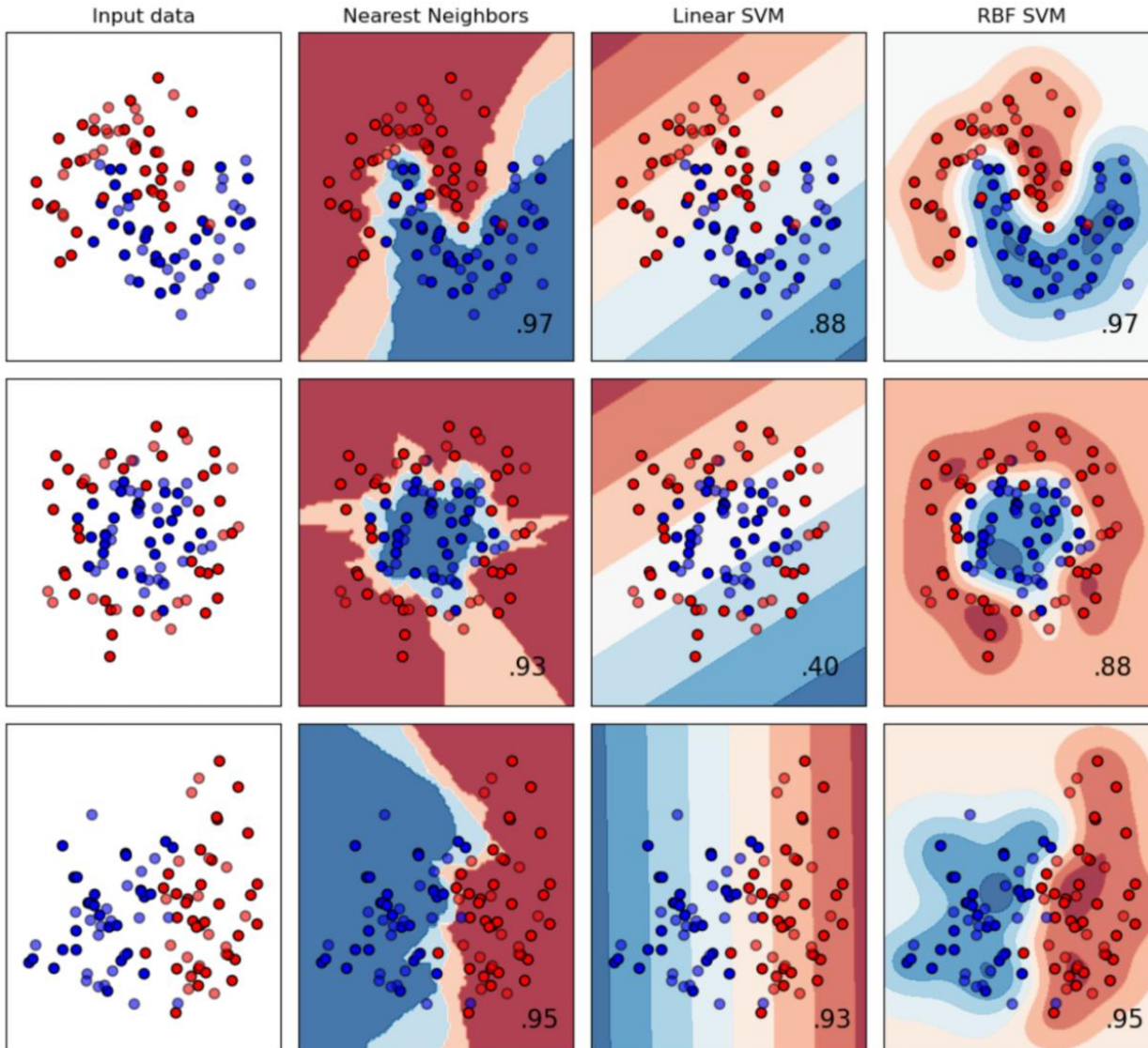
기존의 계산 시간을 극복할 수 있는 “효율성” (대리모델)



다수의 시나리오를 빠르게 탐색할 수도 있음

인공지능에게 기대하는 바

□ 인간이 인식하지 못하는 패턴, 현상, 의미 “포착”



Katto & Ohno (1984)

$$q_c'' = q_{co}'' \left(1 + K \frac{h_f - h_{in}}{h_{fg}} \right)$$

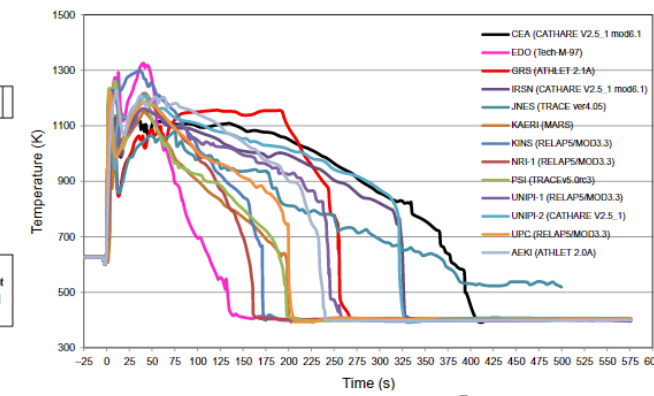
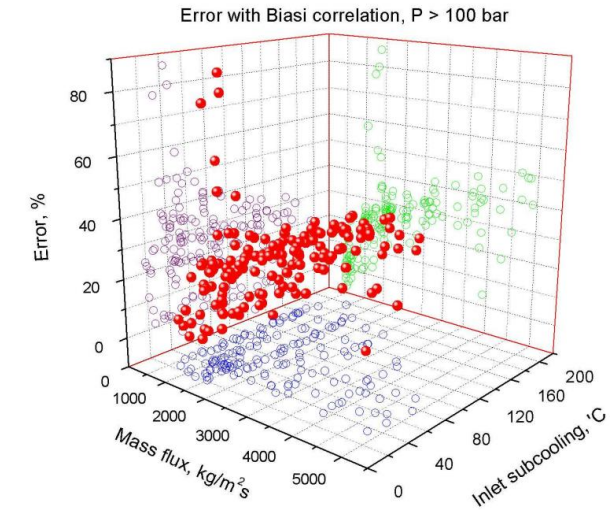
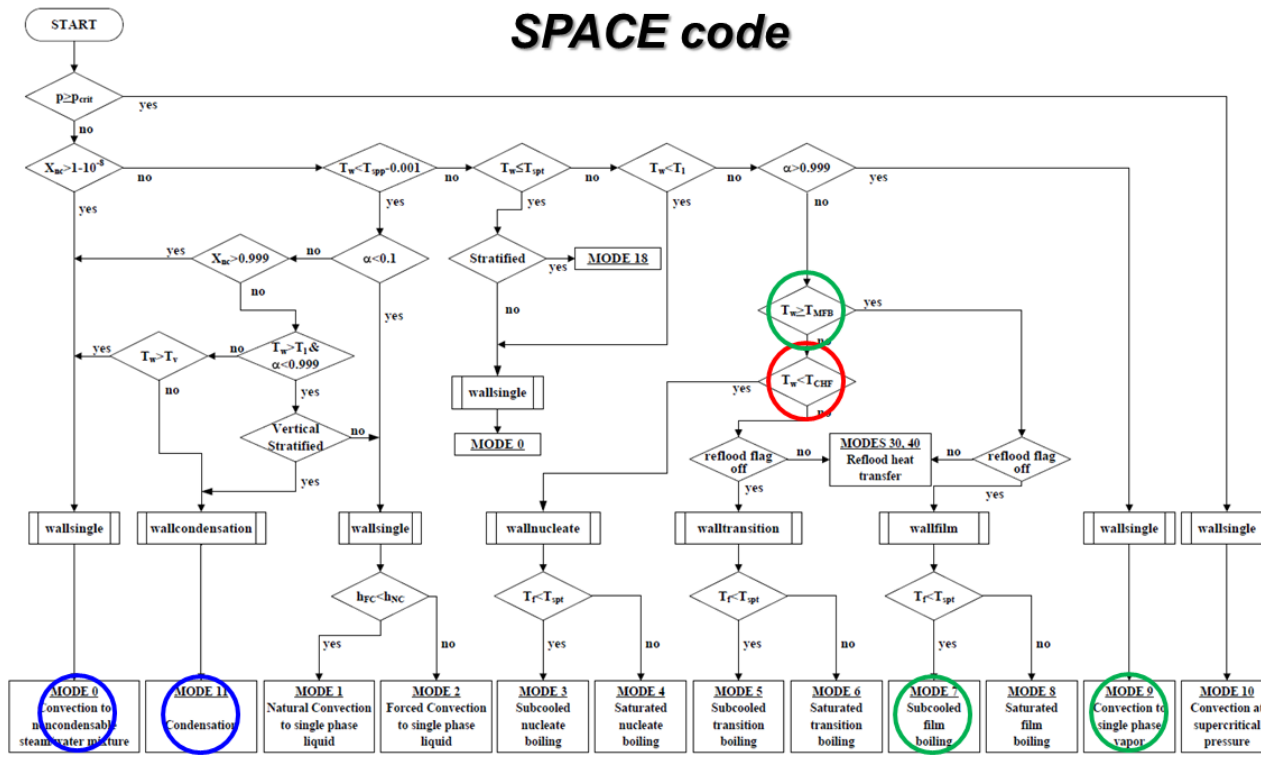
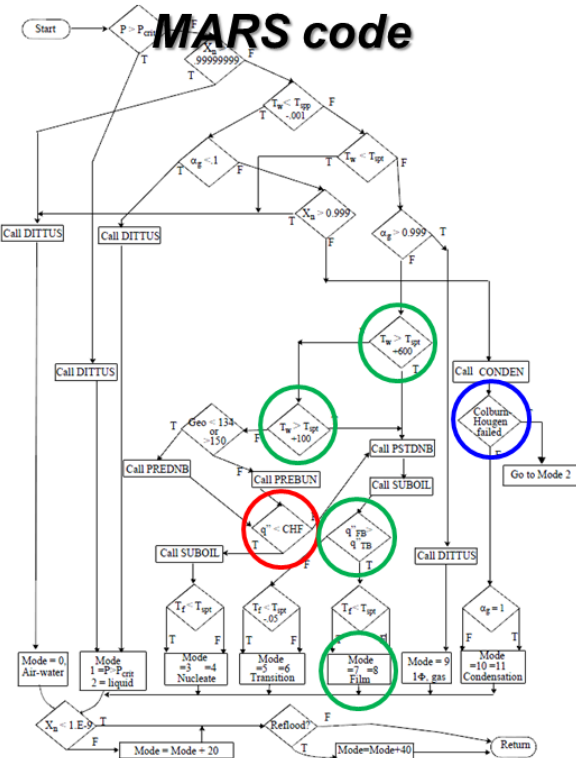
PG-CHF (1986)

$$R = \frac{k_1 F_g}{f(P_r)(dT_r)^{k_2}} \frac{f(P_r, G)f(P_r, X_i)}{f_a} = f(P_r, G, X_i, X, q'', dT_r, F_g)$$

원자로 안전해석 기술 현황

□ 구성 모델(상관식)의 불확실도로 인한 코드 구조 및 예측 거동 상이

- 구성 모델, 열/유동 양식 판별 논리 및 기준이 다름
- 동일 사고 시나리오 및 실험에 대해서 상이한 거동
- “보수적 접근”

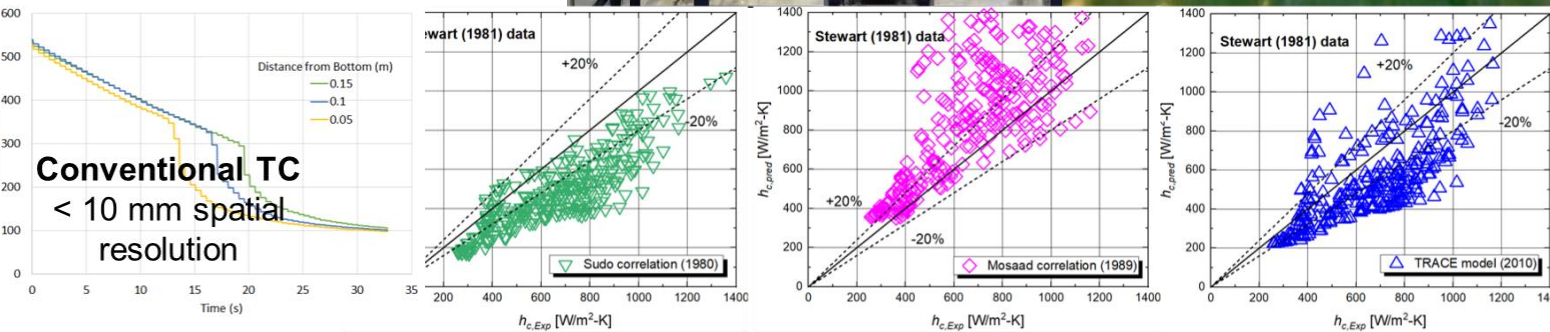
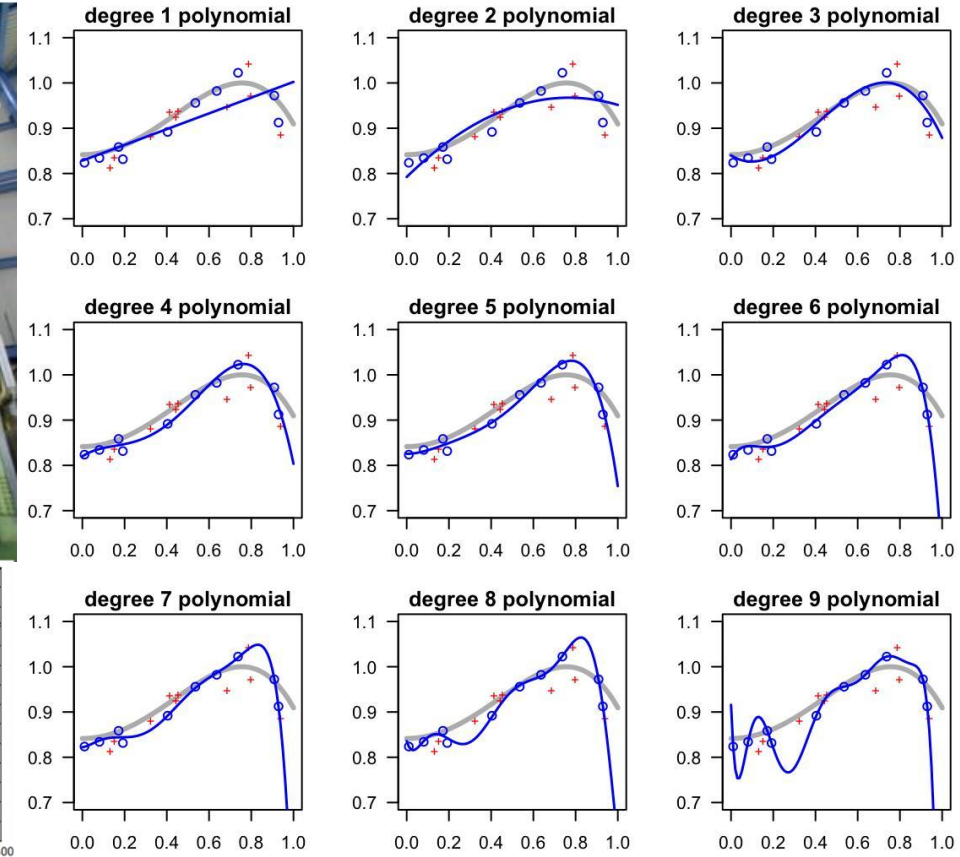
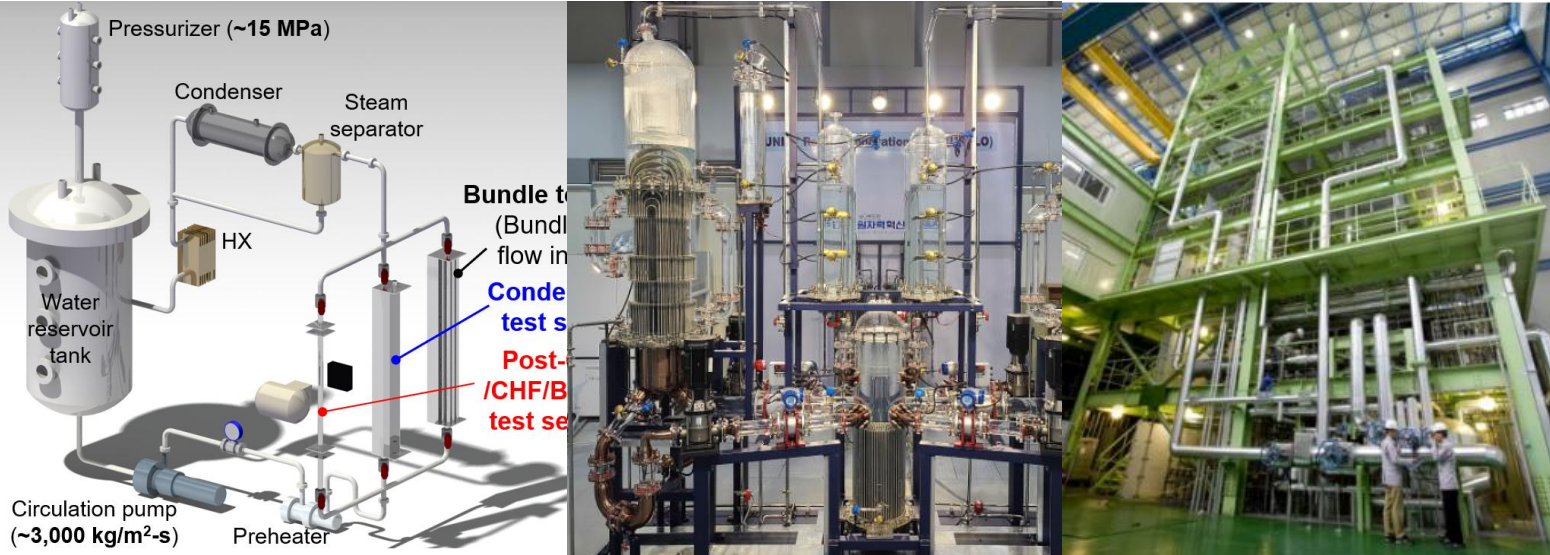


원자로 안전해석 기술 현황

□ 구성 모델(상관식)을 개발하기 위한 연구

열수력 실험 (SET, IET)

실험데이터 회귀



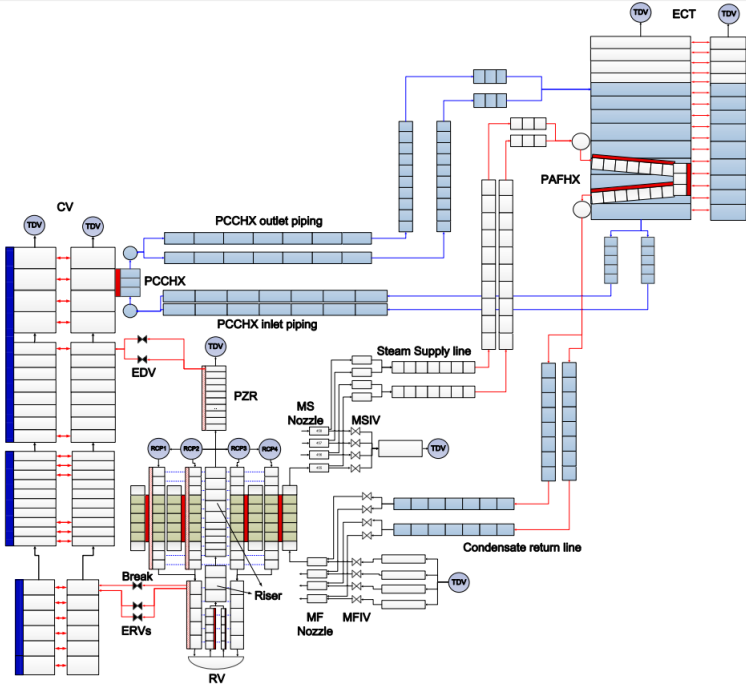
실험 데이터 한계 (개수, 제한적 계측 데이터) 및 불확실성

제한적 회귀 모형

원자로 안전해석 기술 현황

□ 구성 모델(상관식)의 한계에 따른 추가 연구

1D System codes

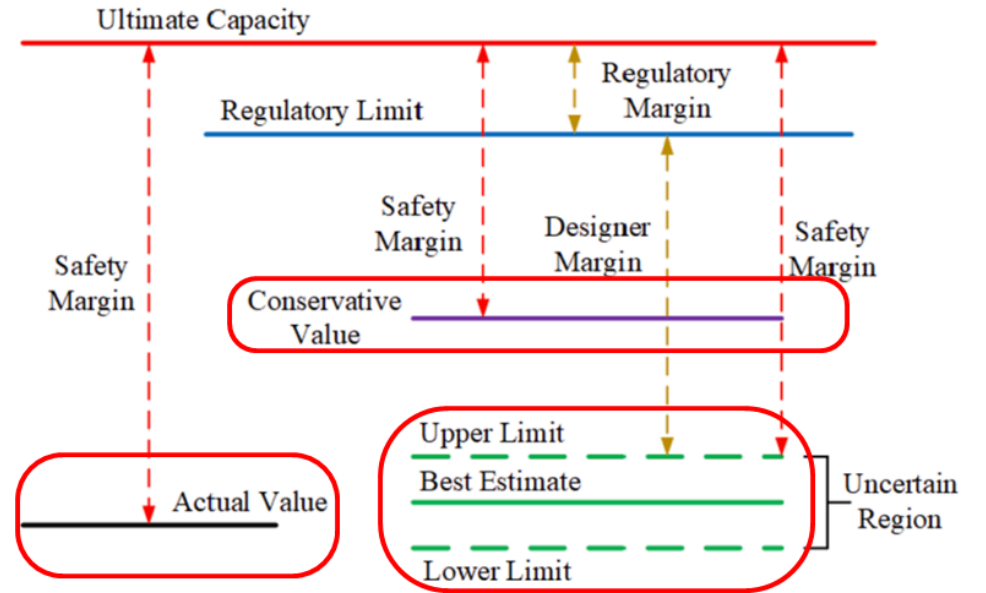


Empirical correlation
→ Limited Predictability

상관식 개선 및 신규 개발
+ 예측 성능 및 불확도 재평가

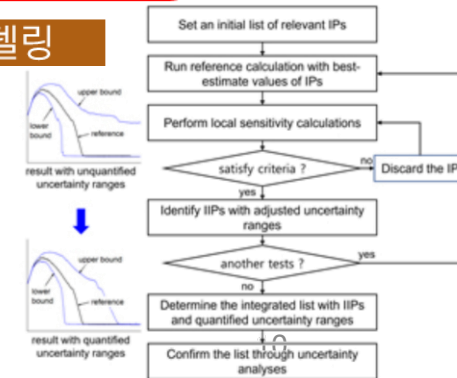
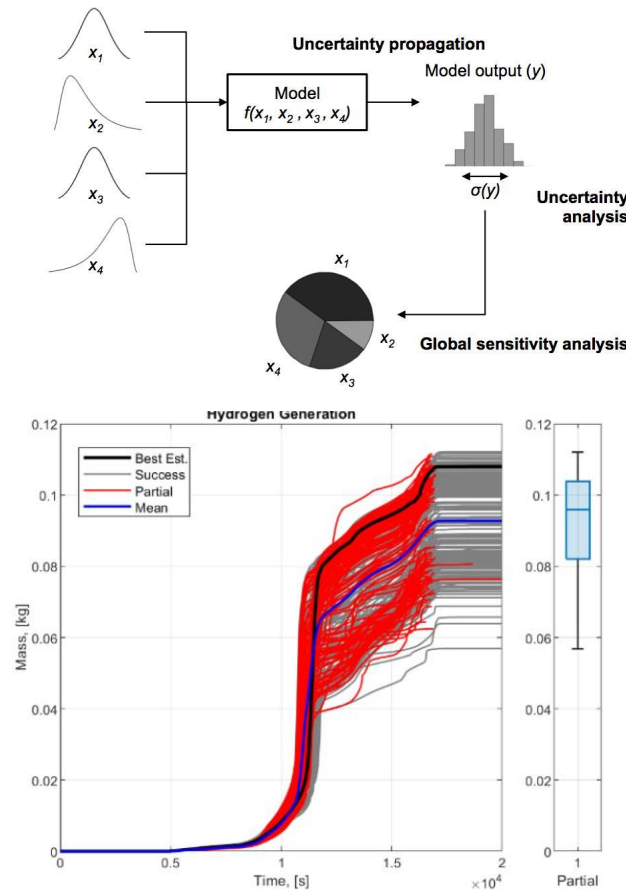
불확실도 평가 (BEPU 등)

안전여유도 및 허용기준 수립 방법론



실험데이터

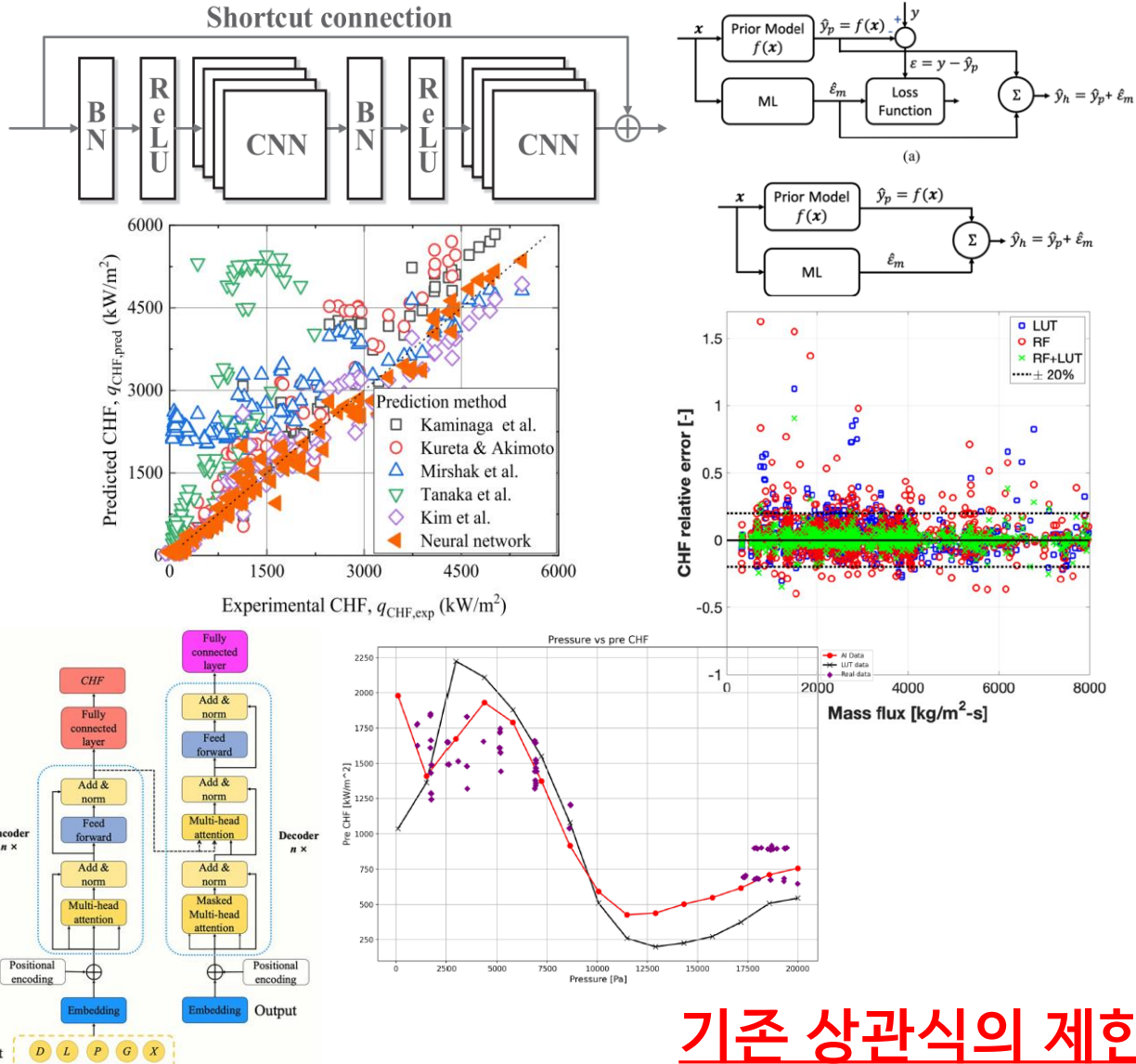
모델링



추가적 연구 제원 소요 및 “보수성”

원자로 안전해석 보조 AI

□ 구성 모델(상관식)을 AI 모델로 대체



NEA Working Paper: Benchmark on Artificial Intelligence and Machine Learning for Scientific Computing in Nuclear Engineering. Phase 1: Critical Heat Flux Exercise Specifications.

Authors: Jean-Marie LE CORRE, Gregory DELPEL, Xu WU, Xingang ZHAO.

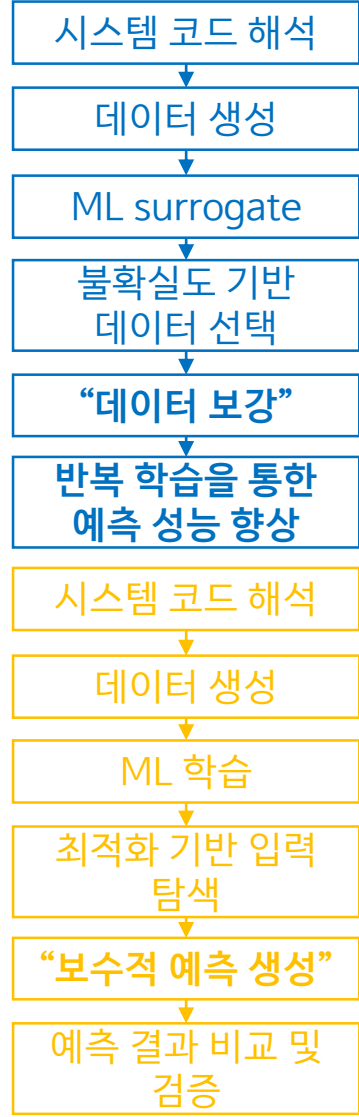
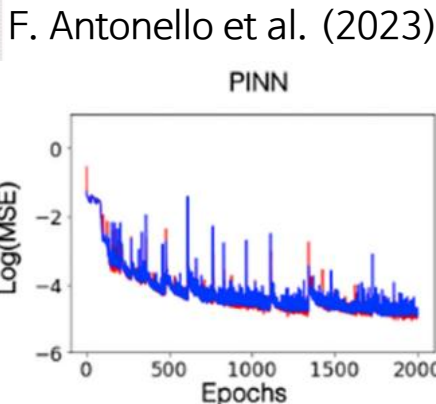
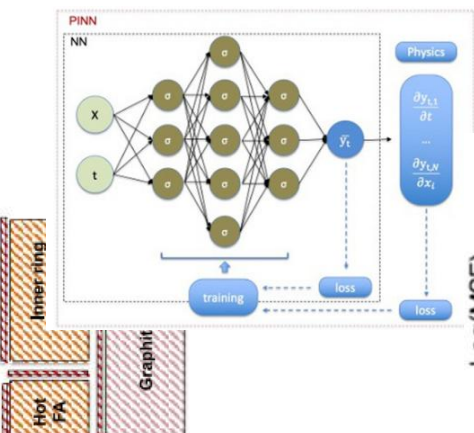
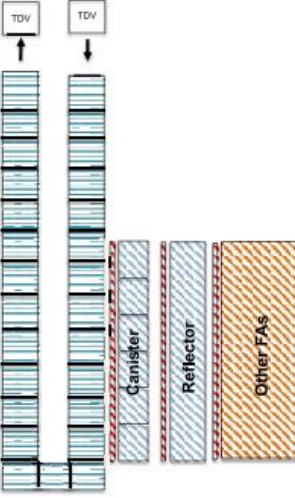
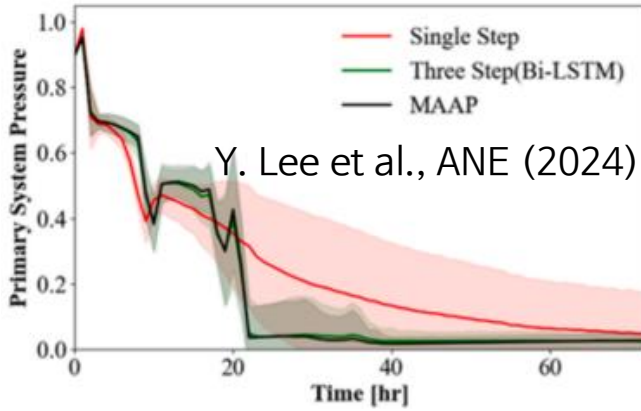
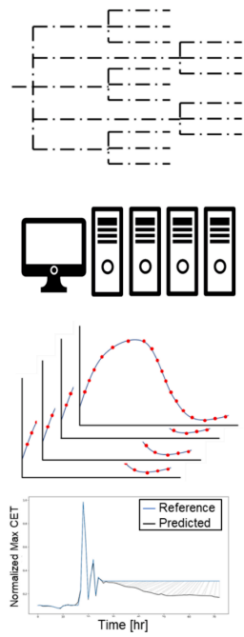
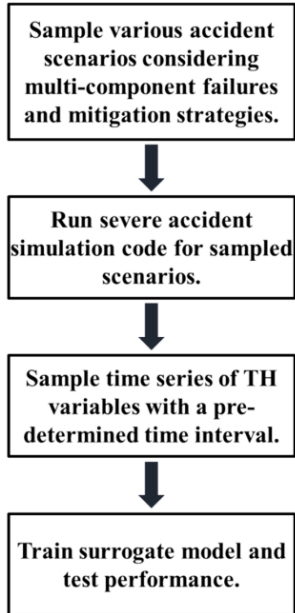
Key attributes: Accuracy, Transparency, Robustness, Trustworthiness.

구성모델(현상)	2021년	2022년	2023년	2024년	2025년	합계
CHF/DNB	8	10	9	14	25	66
2상유동 압력강하	6	8	10	12	12	48
막비등 관련	6	6	7	7	7	33
Net vapor generation	3	5	5	5	5	23
기포울	6	8	9	10	12	45
비등열전달	10	12	12	16	16	66
총합	39	49	52	64	77	281

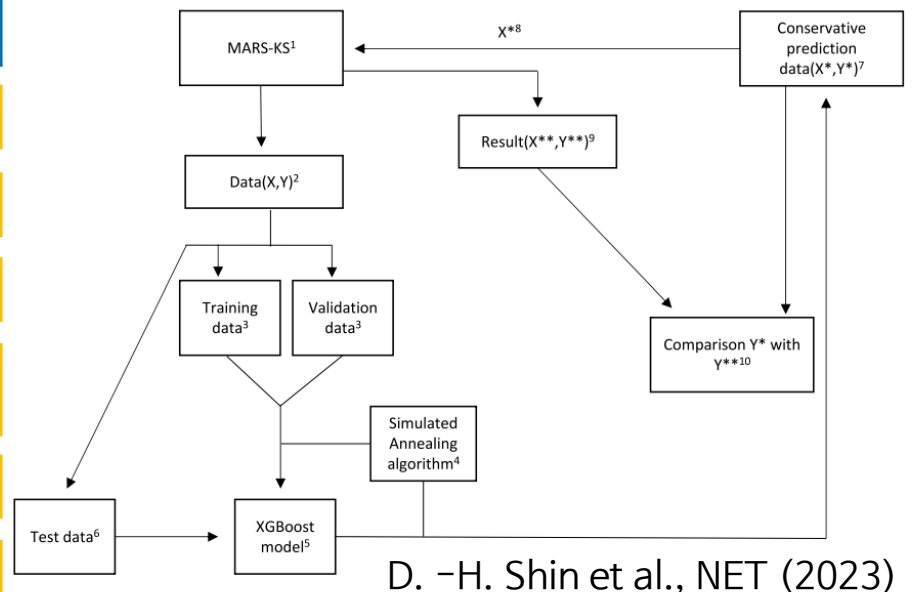
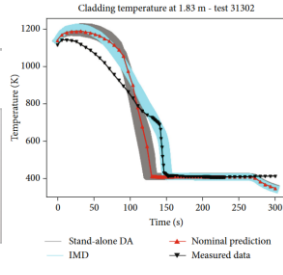
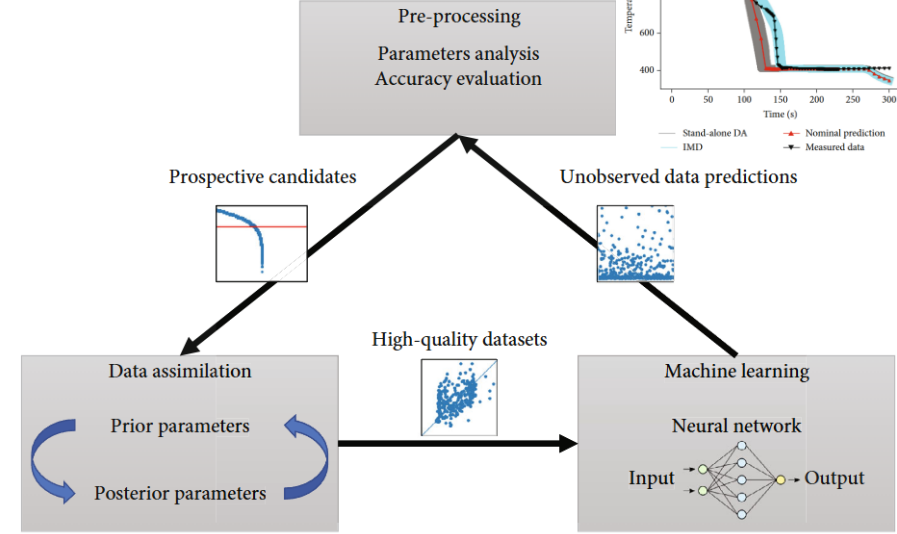
기존 상관식의 제한적 회귀 성능을 극복

원자로 안전해석 보조 AI

시스템 해석 코드 및 구성모델 Surrogate



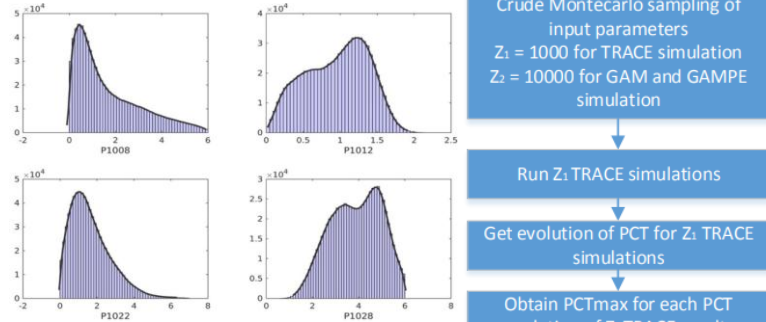
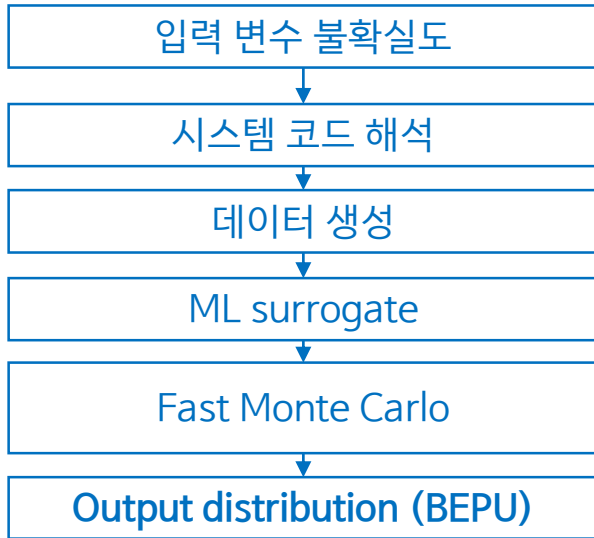
N. H. Tiep et al., IJER (2024)



D. -H. Shin et al., NET (2023)

원자로 안전해석 보조 AI

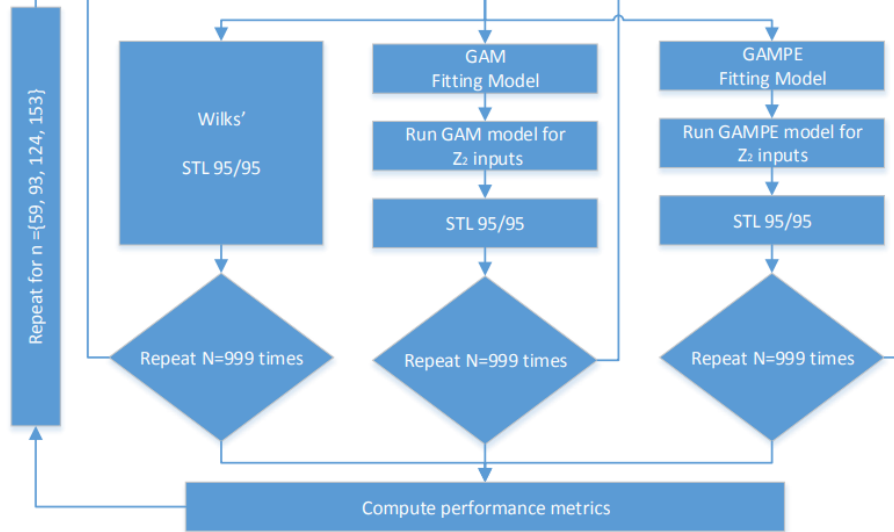
시스템 해석 코드 기반 신속 BEPU를 위한 대리 모델



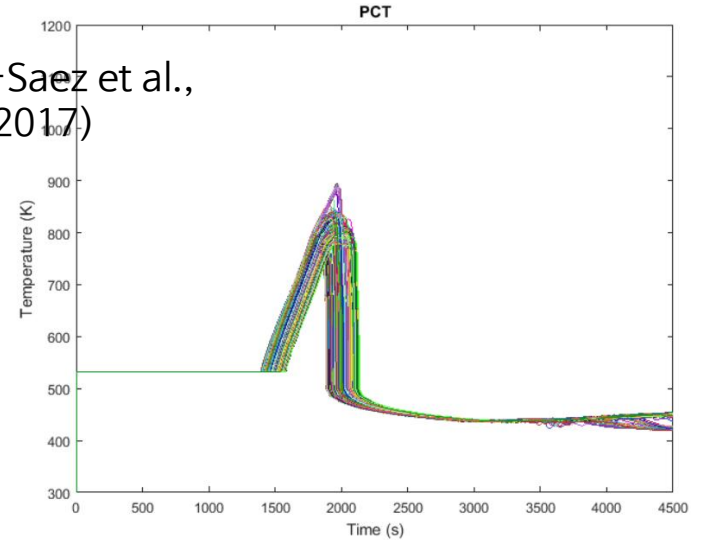
Crude Montecarlo sampling of input parameters
 $Z_1 = 1000$ for TRACE simulation
 $Z_2 = 10000$ for GAM and GAMPE simulation

Run Z_1 TRACE simulations
 Get evolution of PCT for Z_1 TRACE simulations
 Obtain PCTmax for each PCT evolution of Z_1 TRACE results

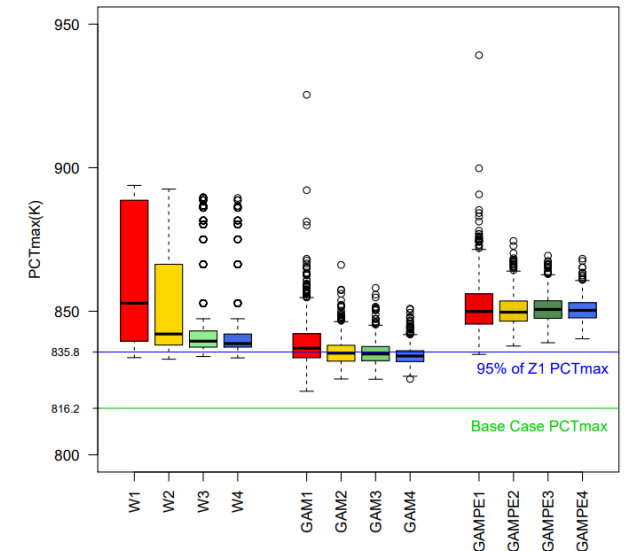
Random sampling of n input and output (PCTmax) values from Z_1



F. Sanchez-Saez et al.,
 NED(2017)



Wilks' and GAMs - STL 95/95

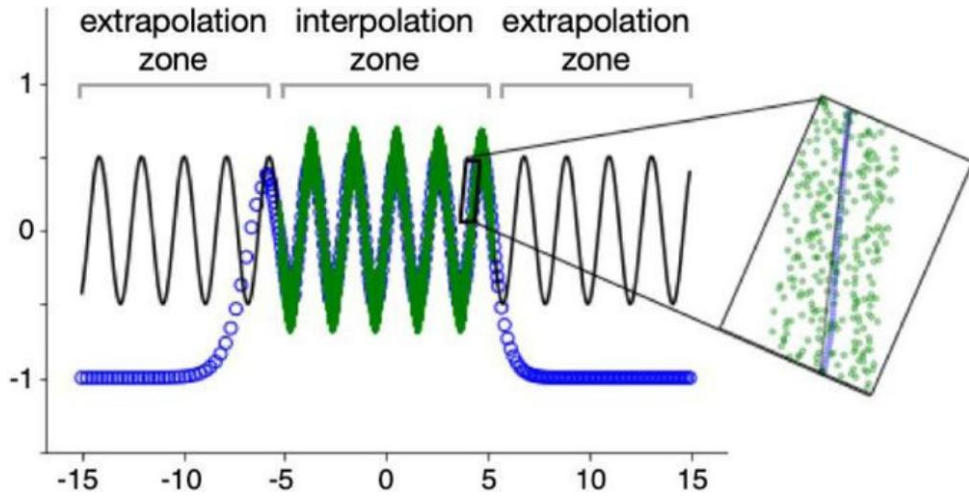


빠른 탐색 → 빠른 BEPU 수행

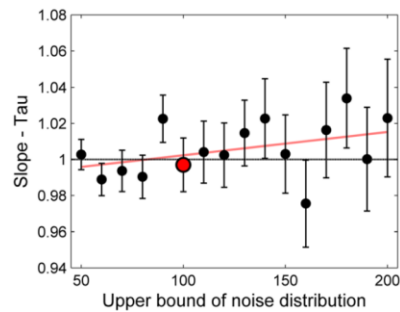
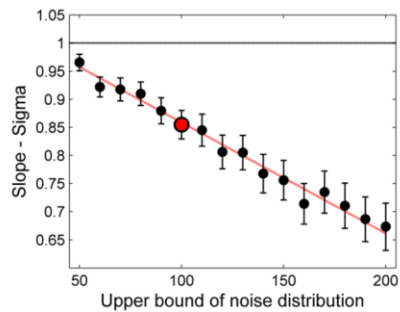
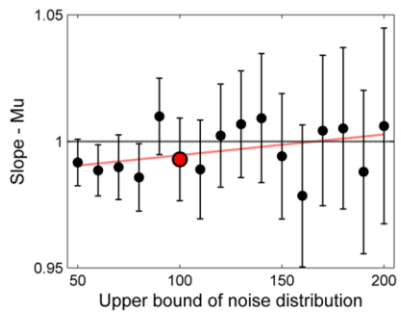
기존 원자로 안전해석 보조 SI 한계

신뢰도를 평가하기 어려움 & 변화에 취약

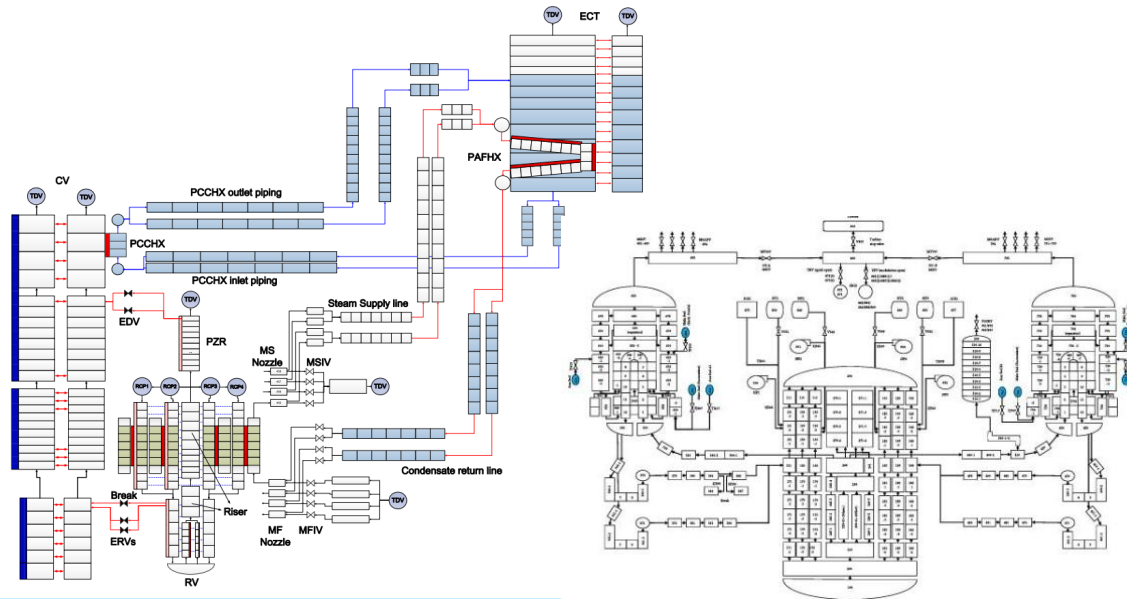
일반화 (외삽) 성능



불확실도 (신뢰도) 정량화



시계열 데이터 및 대표 변수 데이터 기반탐색



시계열 데이터 추출 위치, 변수 조합 등에 따른 민감도 → 최적 조합 탐색 필요

공간에 대한 분해능 없음 (HL, CL, PRZ 등 위치 labeling 필요)

기존 원자로 안전해석 보조 SI 한계

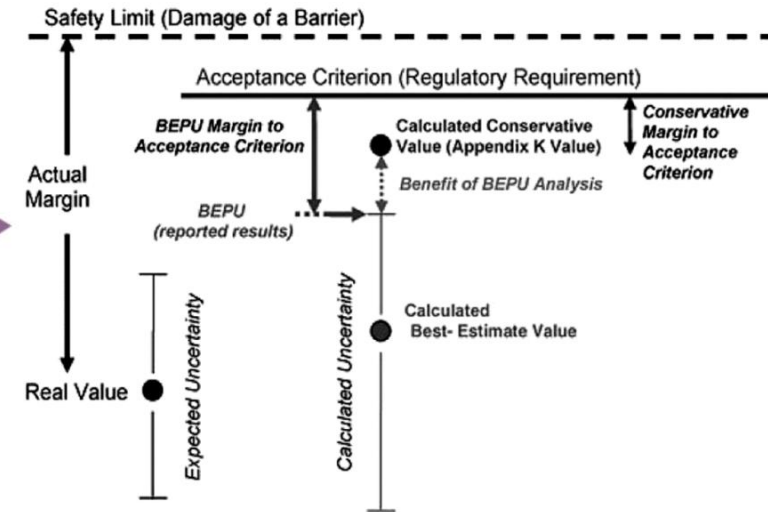
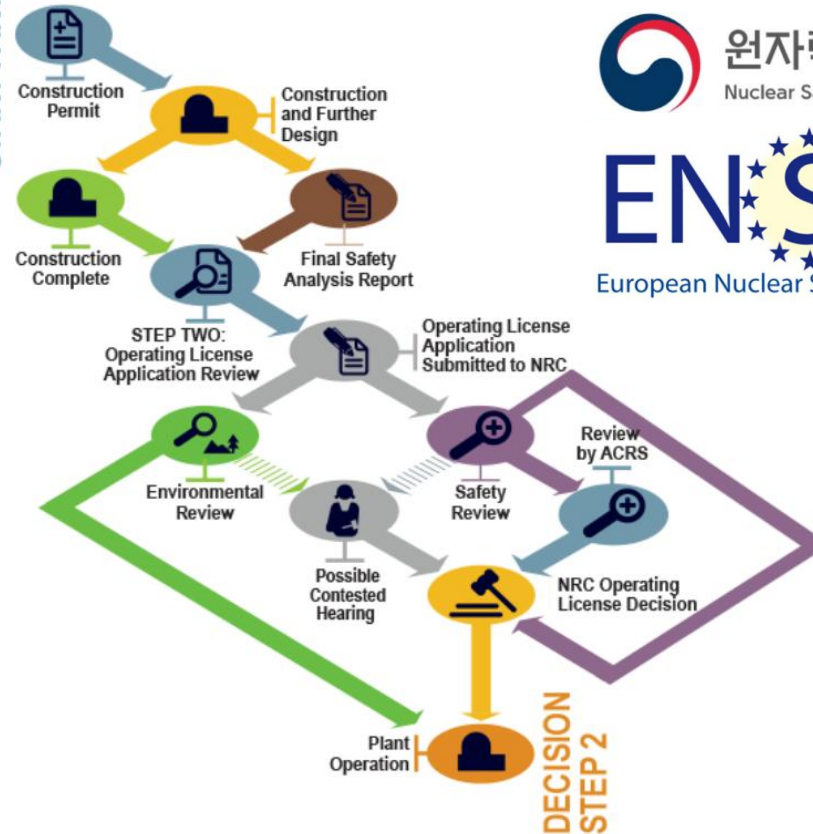
□ (진입 장벽)인공지능 기반 기술 또한 설계 인가 및 상업운전 시 적용을 위해서는, 원자력규제기관의 허가를 취득해야함

원자력 발전소 운영을 위한 인허가 절차

START STEP 1



START PART 2

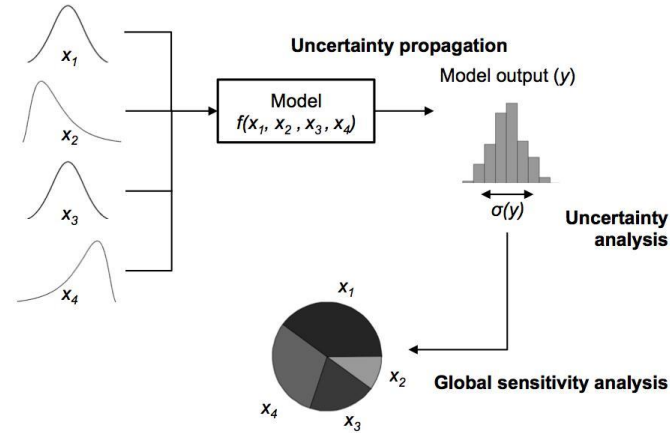
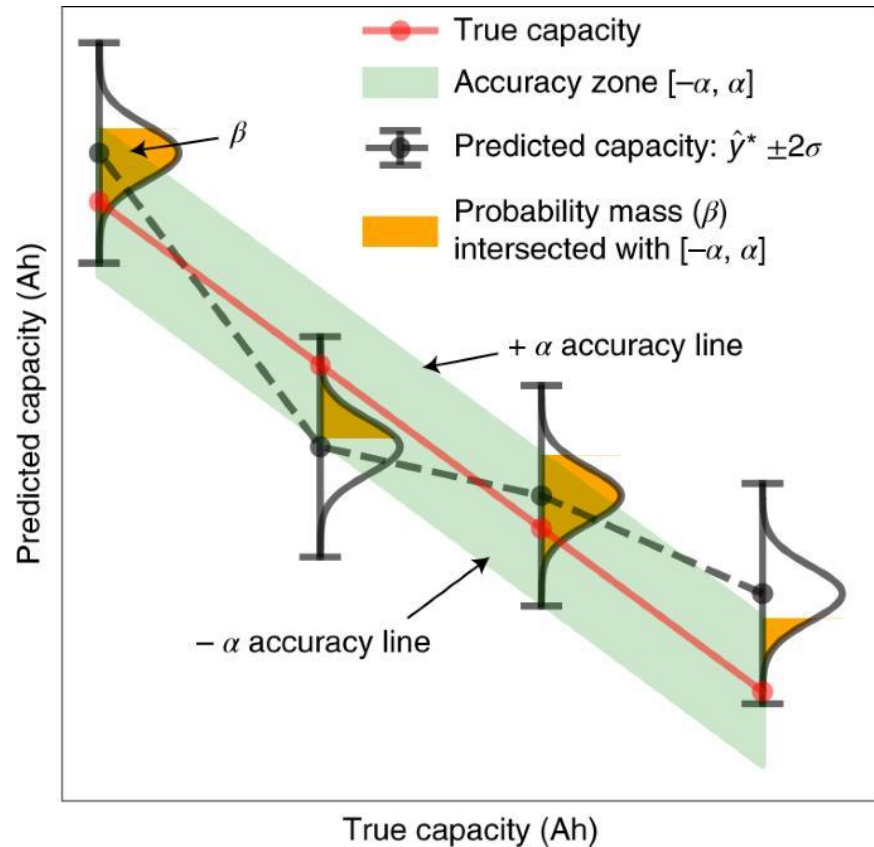


원자로 안전해석 보조 AI 한계 극복 방법

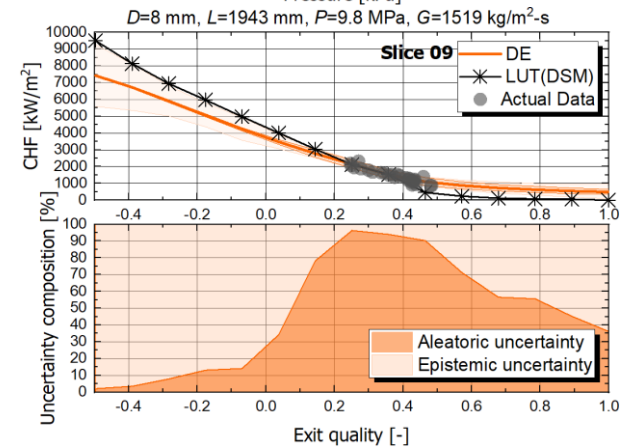
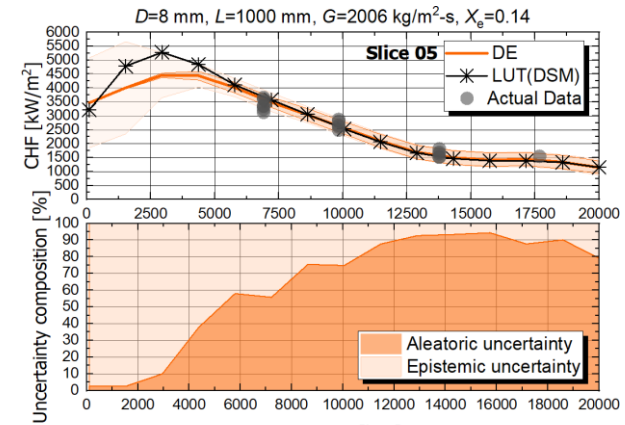
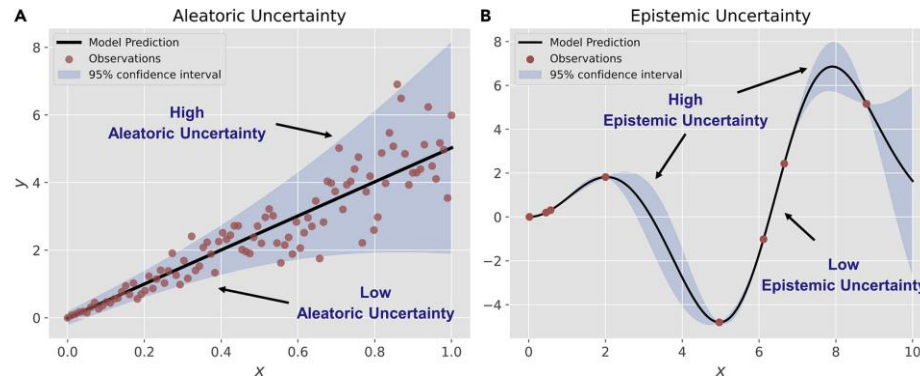
□ AI 모델 예측에 대한 불확실도 정보 제공 & 불확실도의 신뢰도 평가

예측에 대한 불확실도 (신뢰구간) 정량화

예측 및 불확실도에 대한 인자별 영향성



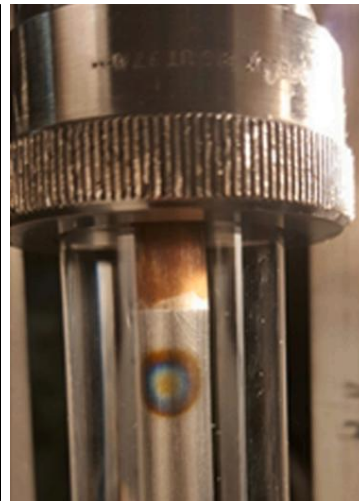
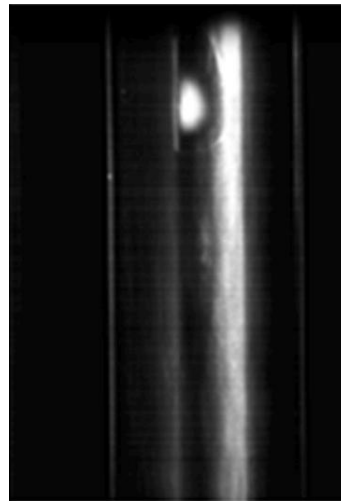
불확실도 원인 분석



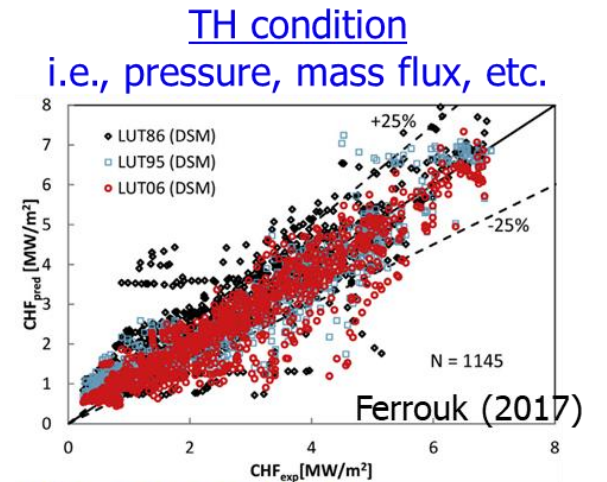
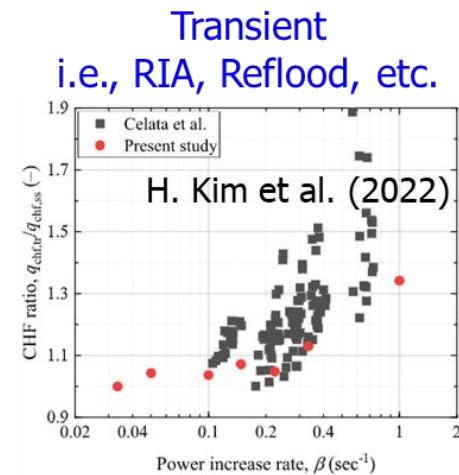
불확실도 측면 Trustworthy 구성 모델 대체 AI 모델 개발

□ 불확도 정량화 가능 AI 모델 - 임계열유속(CHF) case

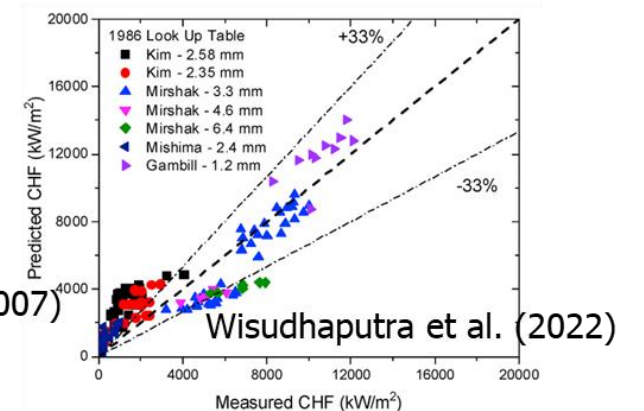
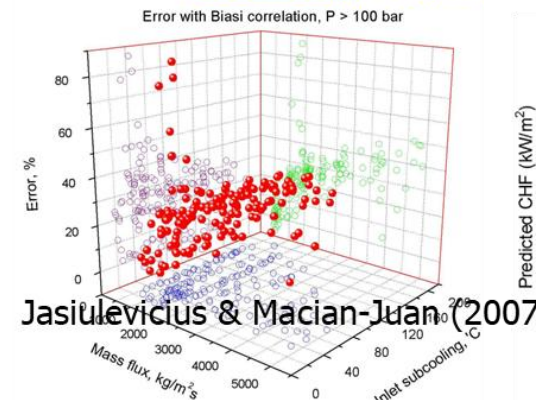
- 임계열유속은 경수로의 대표적인 안전지표(FOM)
- 그러나 현재 원자로 시스템 안전해석 코드 내 CHF 예측 불확실도 존재



Moreira et al. (2022)



Geometry, i.e., Bundle, rectangular, etc.

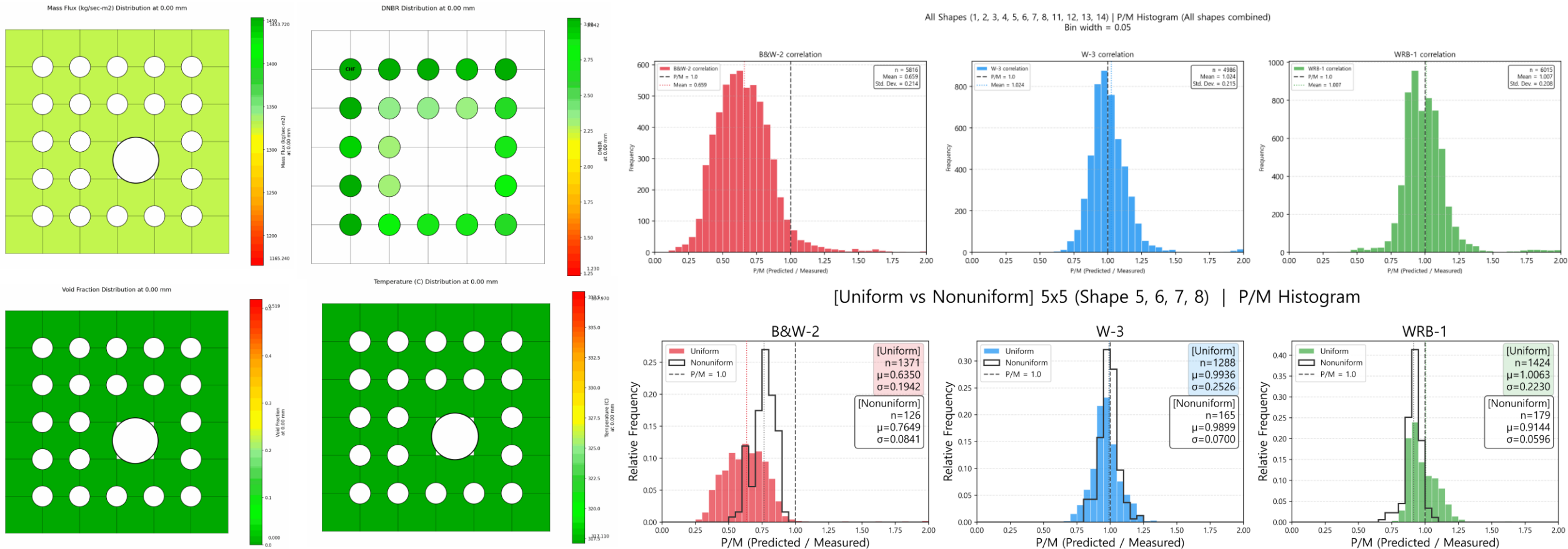


불확실도 측면 Trustworthy 구성 모델 대체 AI 모델 개발

□ 기존 구성 모델(상관식) 및 시스템/부수로 해석 코드 예측 성능 한계

- 현재 원자로 시스템 안전해석 코드 내 CHF 예측 불확실도 존재

→ 보수적인 평가 방법론



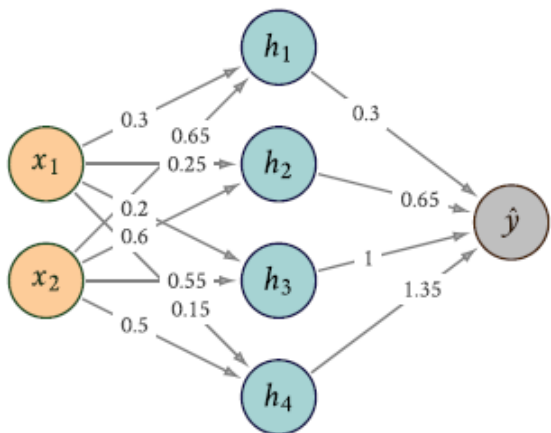
불확실도 측면 Trustworthy 구성 모델 대체 AI 모델 개발

□ 불확도 정량화 가능 AI 모델 - 임계열유속(CHF) case

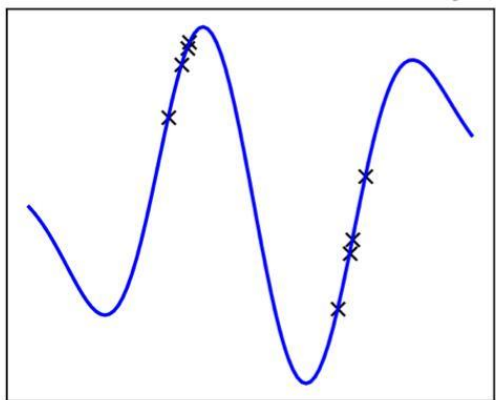
- 확률론적 모델(BNN, MCD, DE, GP 등) 기반 CHF 예측 프레임워크 개발

→ 예측 + 불확실도

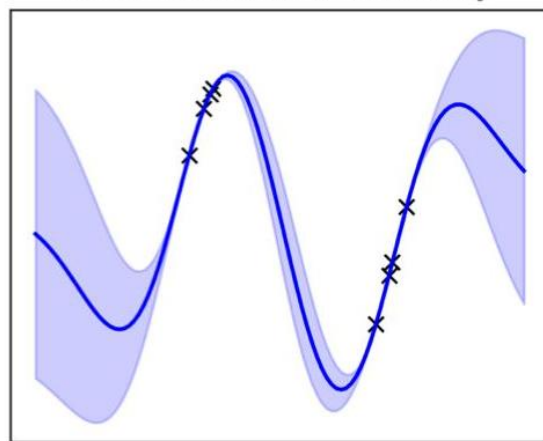
Deterministic NN



Prediction without uncertainty

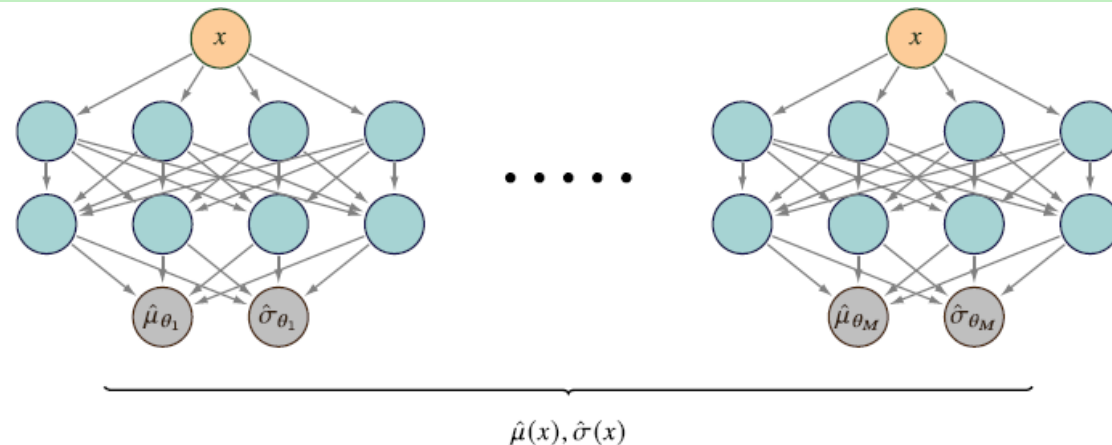


Prediction with uncertainty

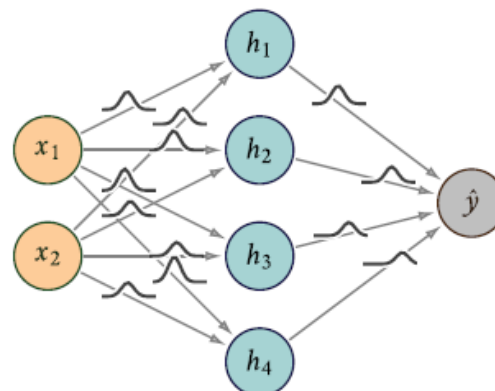


VS

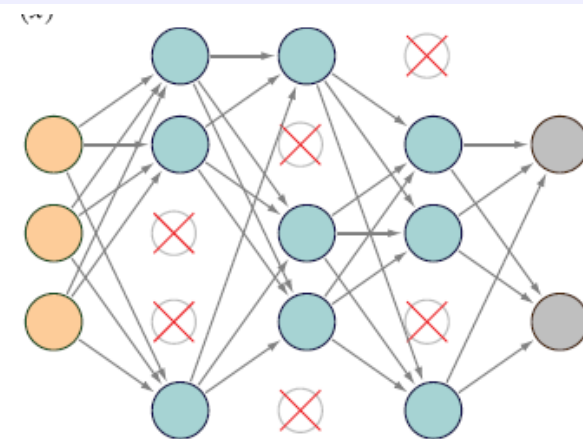
Deep Ensemble (Heterogeneous)



Bayesian Neural Network (BNN)



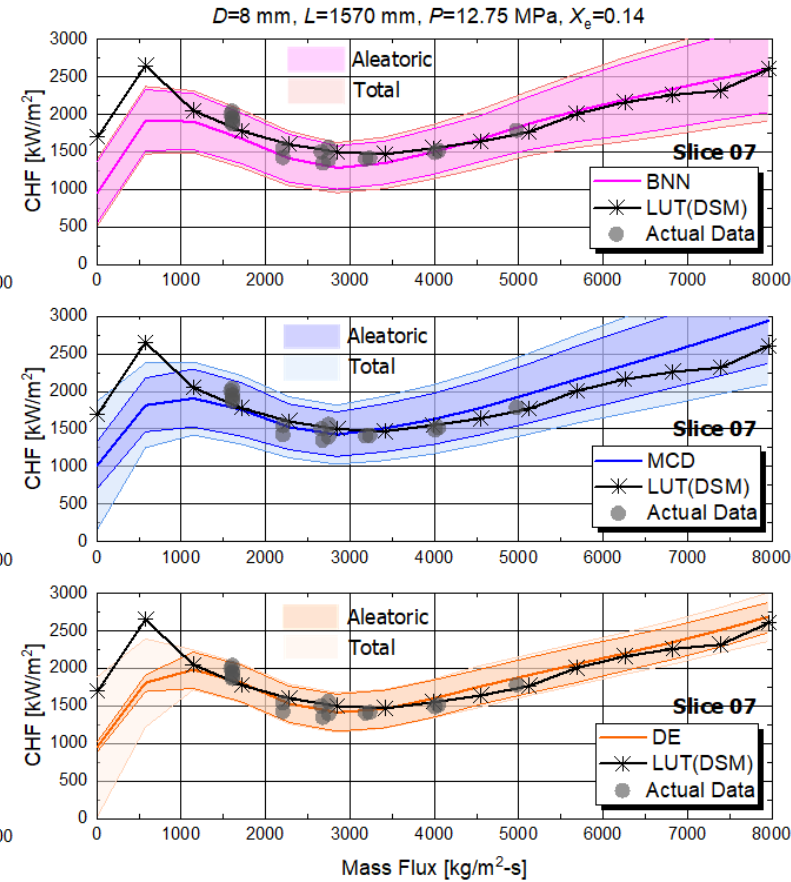
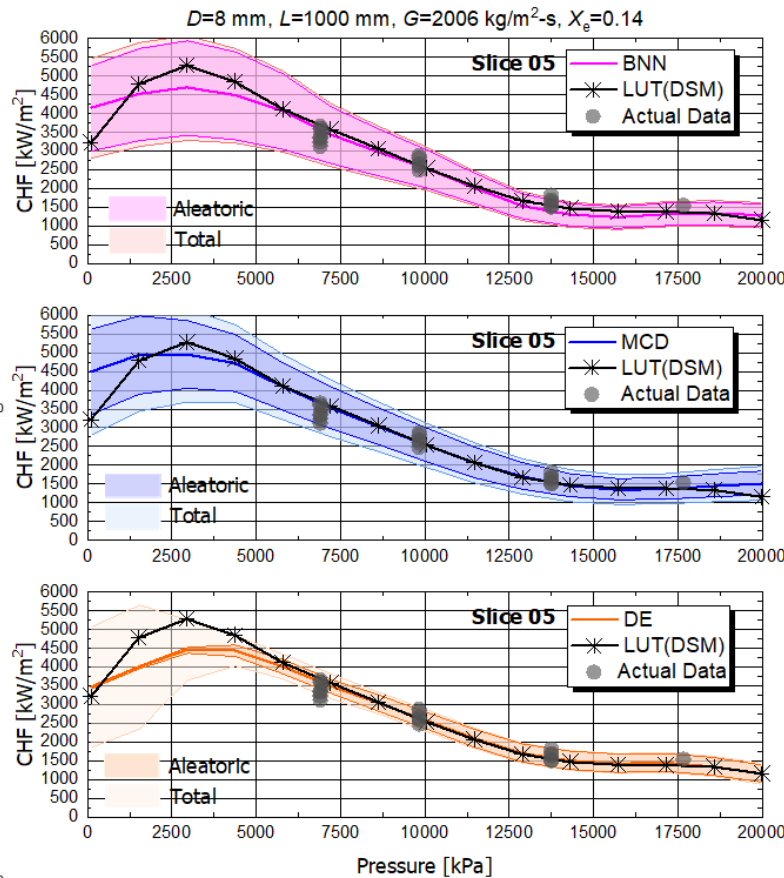
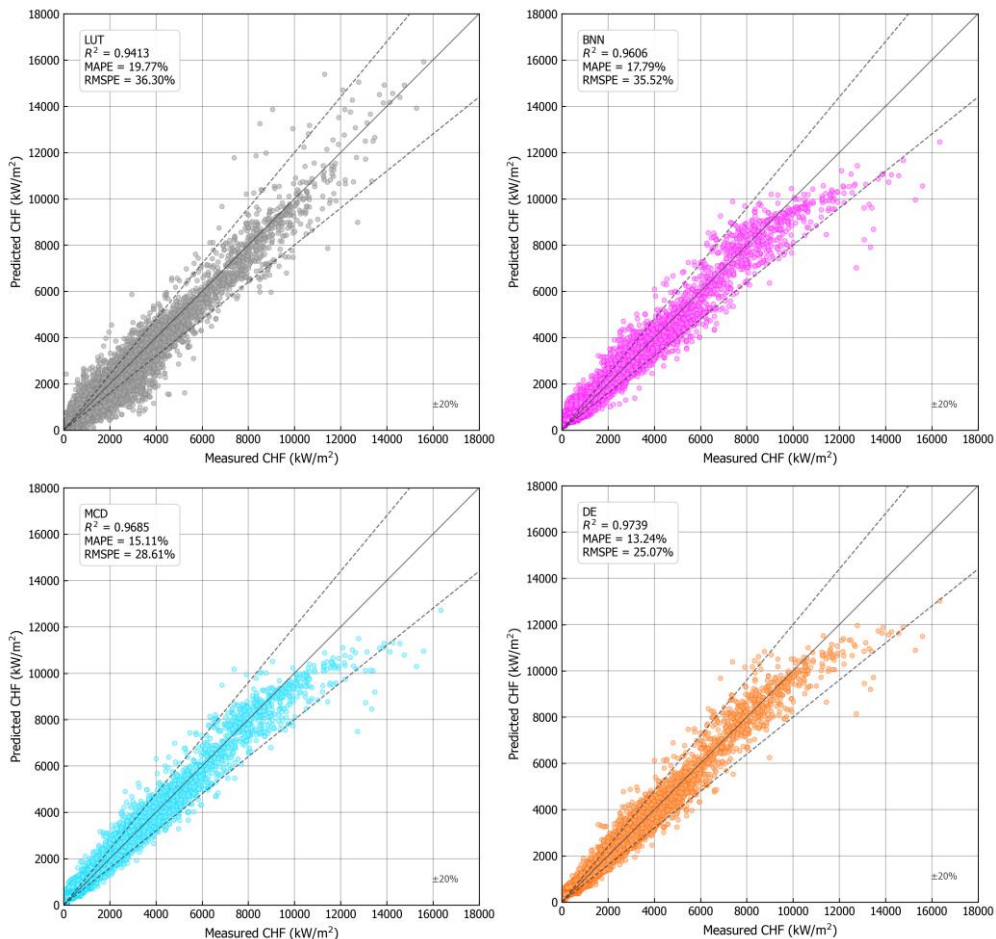
Monte Carlo Dropout (MCD)



불확실도 측면 Trustworthy 구성 모델 대체 AI 모델 개발

□ 임계열유속 (CHF) case - 회귀 성능 및 불확실도 정보 제공

- 기존 상관식 대비 우수한 예측 성능 + Deterministic NN 대비 불확실도 정보 제공 → 그러나, 모델이 제공하는 불확실도 (신뢰구간) 정보를 신뢰 가능한가?



불확실도 측면 Trustworthy 구성 모델 대체 AI 모델 개발

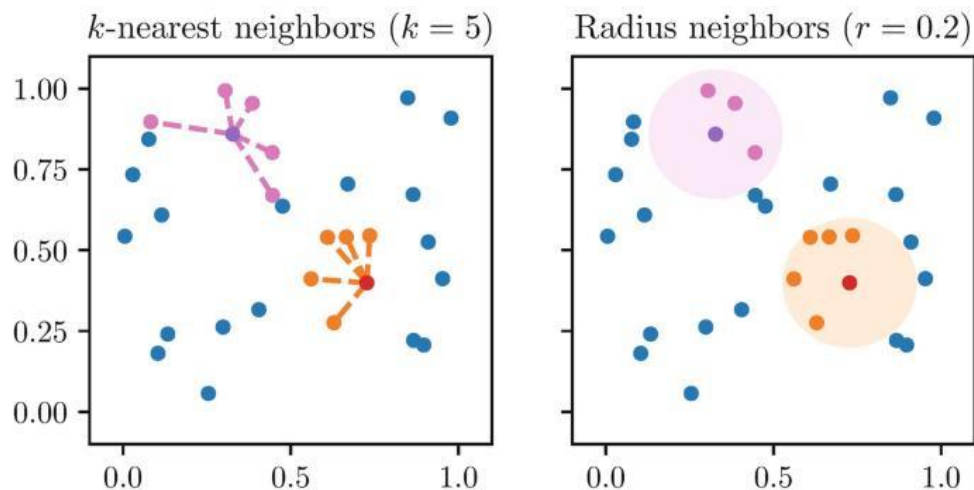
□ 불확실도 (신뢰구간) 정보 신뢰 평가 방법 (외삽 영역)

- Epistemic uncertainty는 데이터 수와 반비례해야 함.

→ KNN_{100} 거리, Mahalanobis 거리, N_{nearby} 와 Epistemic uncertainty 비교

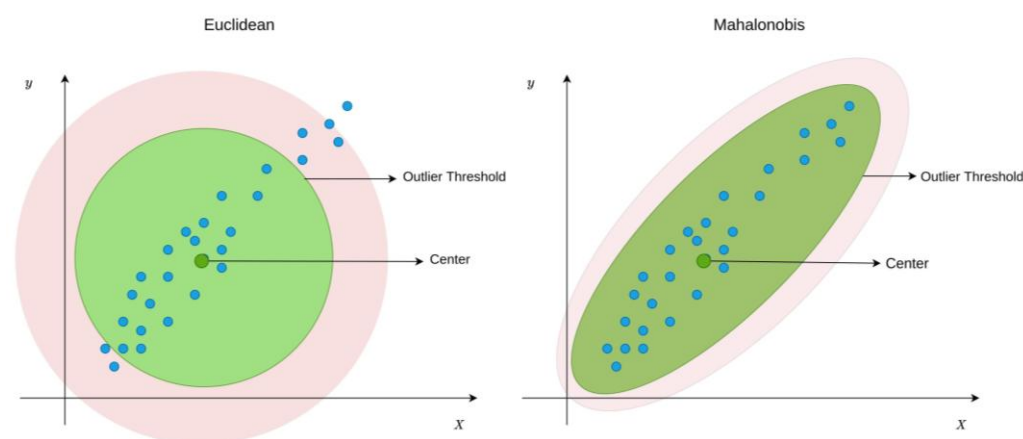
주변 100개 데이터까지의 KNN 거리 (Local)

$$d_{KNN}(x) = \frac{1}{K} \sum_{k=1}^K |x - x_{(k)}|_2$$



Mahalanobis 거리 (global)

$$d_{Mahal}(x) = \sqrt{(x - \mu)^T \Sigma^{-1} (x - \mu)}$$



KNN/Mahalanobis와 epistemic 불확실도간 상관관계 (Spearman rank correlation coefficient)

$$\rho_{KNN} = \rho_s(\sigma^2, d_{KNN})$$

$$\rho_{Mahal} = \rho_s(\sigma^2, d_{Mahal})$$

불확실도 측면 Trustworthy 구성 모델 대체 AI 모델 개발

□ 불확실도 (신뢰구간) 정보 신뢰 평가 방법 (외삽 영역)

- KNN_{100} 거리, Mahalanobis 거리, N_{nearby} 와 Epistemic uncertainty 비교
- $N_{nearby}(r=1.5)$ 가 증가함에도 Epistemic 감소, d_{KNN100} 이 감소함에도 변화 미비

주변 100개 데이터까지의 KNN 거리

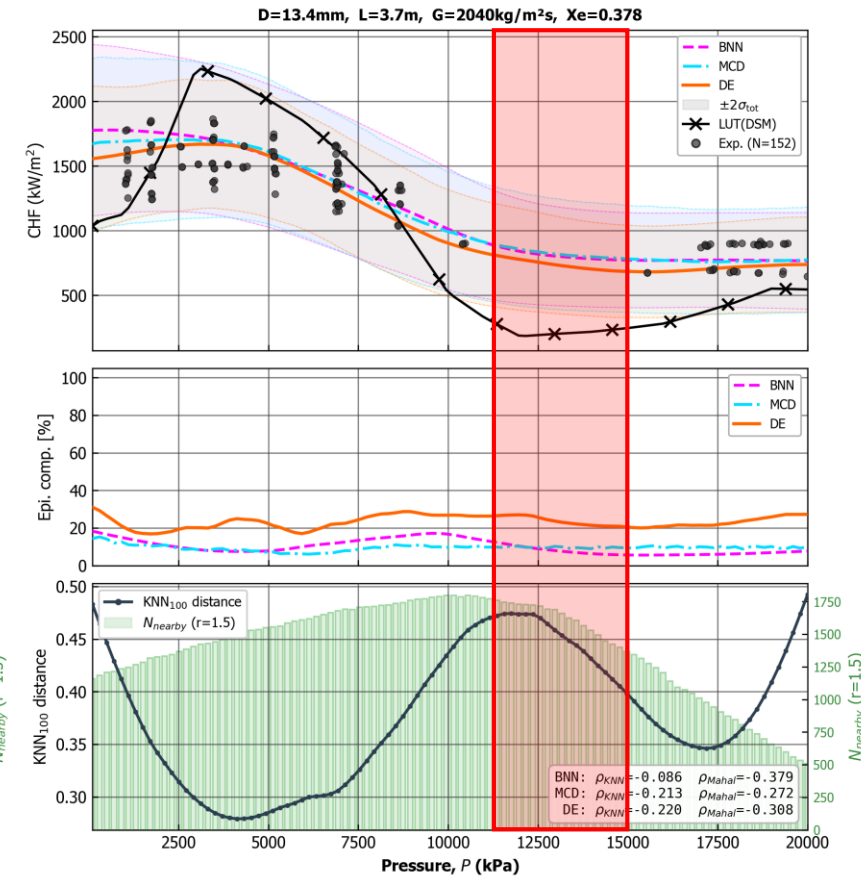
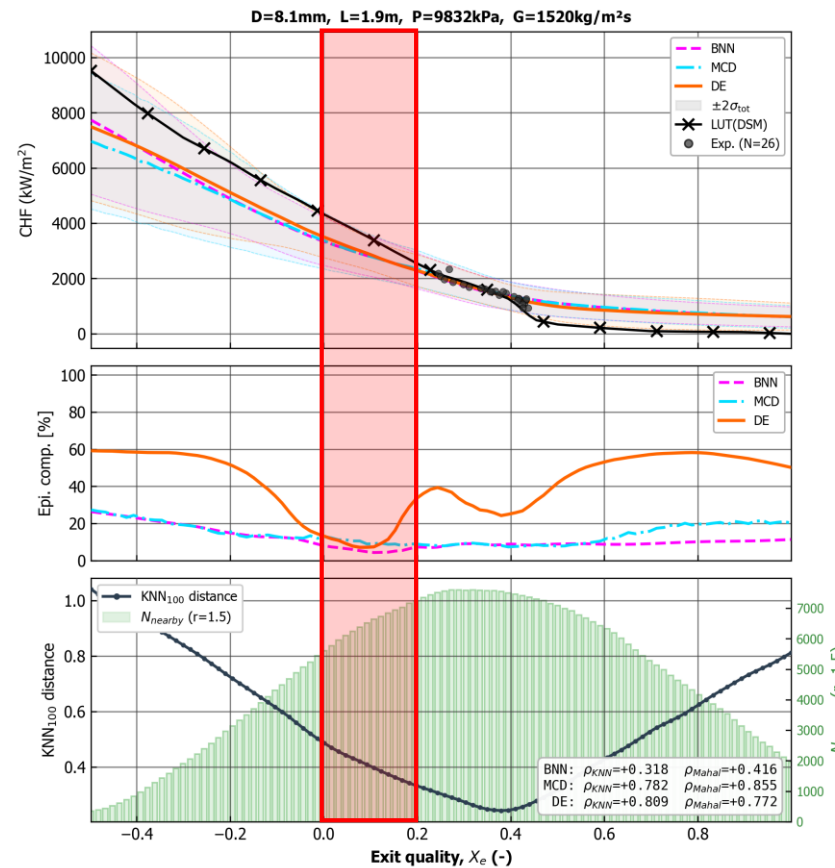
$$d_{KNN}(x) = \frac{1}{K} \sum_{k=1}^K |x - x_{(k)}|_2$$

Mahalanobis 거리

$$d_{Mahal}(x) = \sqrt{(x - \mu)^T \Sigma^{-1} (x - \mu)}$$

KNN와 epistemic 불확실도간 상관관계
(Spearman rank correlation coefficient, Pearson correlation)

$$\rho_{KNN} = \rho_s(\sigma^2, d_{KNN})$$



불확실도 측면 Trustworthy 구성 모델 대체 AI 모델 개발

□ 불확실도 (신뢰구간) 정보 신뢰 평가 방법 (내삽 영역)

- AUSE와 ECE
- 실제 error와 uncertainty간 관계성 평가

Area Under the Sparsification Error (AUSE)

$$AUSE(S) = \int_0^1 \frac{MAE(S_V^U(\alpha))}{MAE(S)} - \frac{MAE(S_V^{AE}(\alpha))}{MAE(S)} d\alpha$$

“Uncertainty가 높은” 데이터 제거 시 데이터 subset에 대한 MAE
 “실제 error가 큰” 데이터 제거 시 데이터 subset에 대한 MAE

Uncertainty가 높은 데이터부터 제거
 → 남은 데이터 error 계산
 → AUSE 낮을수록 uncertainty 신뢰 가능
 ※ 순서 기반 → Epistemic 검증

Expected Calibration Error (ECE)

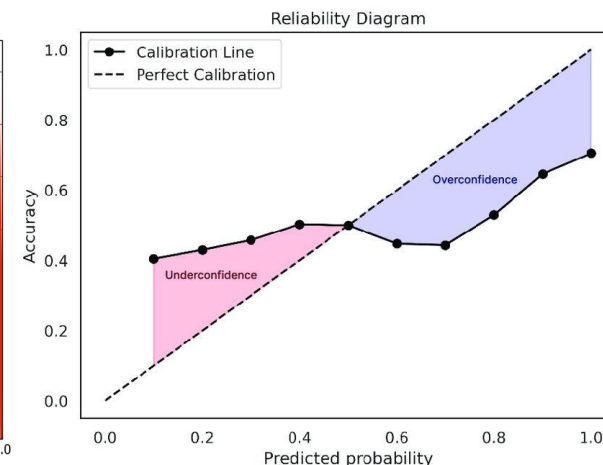
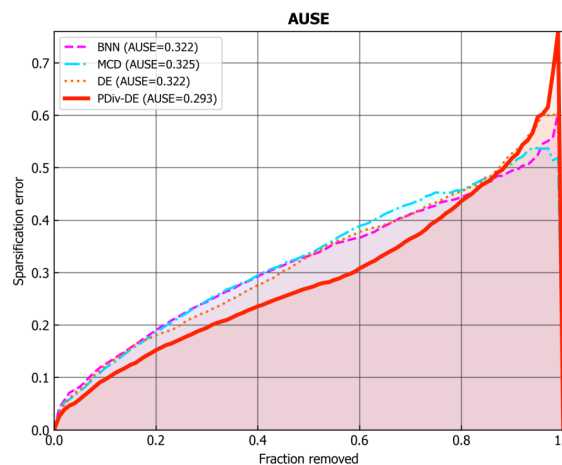
$$ECE = \sum_{k=1}^K \frac{|B_k|}{N} |err(B_k) - unc(B_k)|$$

$$err(B_k) = \frac{1}{|B_k|} \sum_{i \in B_k} |y_i - \mu(x_i)|$$

실제 error 평균

$$unc(B_k) = \frac{1}{|B_k|} \sum_{i \in B_k} \sigma(x_i)$$

예측 uncertainty 평균



불확실도 측면 Trustworthy 구성 모델 대체 AI 모델 개발

□ 일반적 Probabilistic model의 불확실도 정보 신뢰 한계 원인

- 불확실도 자체가 예측 결과에도 비례하는 구조

$$y^{(t)}(x) = f(x; \tilde{W}^{(t)}) \quad \text{Dropout이 적용된 weight}$$

$$\tilde{W} = W \cdot z$$

$$z \sim \text{Bernoulli}(1 - p_D)$$

1차원 가정 시

$$y = wx$$

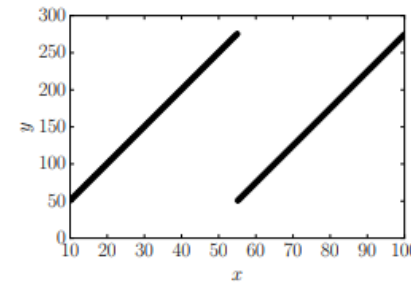
$$y^{(t)} = (wz)x = z \cdot (wx) = z \cdot y$$

기대값 및 분산

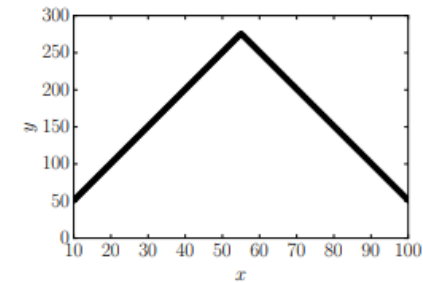
$$\mathbb{E}[y^{(t)}] = \mathbb{E}[z] \cdot y = (1 - p_D)y$$

$$\text{Var}(y^{(t)}) = \text{Var}(z \cdot y) = y^2 \cdot \text{Var}(z) \quad \text{상수 분리}$$

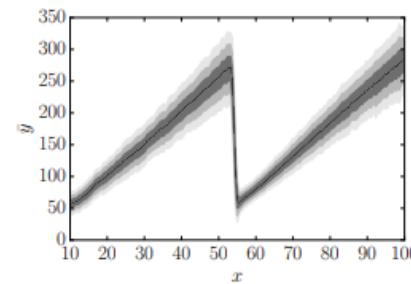
$$\text{Var}(y^{(t)}) = y^2 \cdot p_D(1 - p_D) \quad \because \text{Var}(z) = p_D(1 - p_D)$$



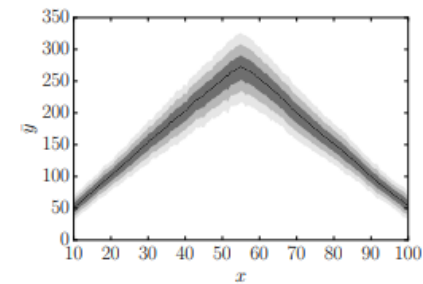
(a) Ground-truth function



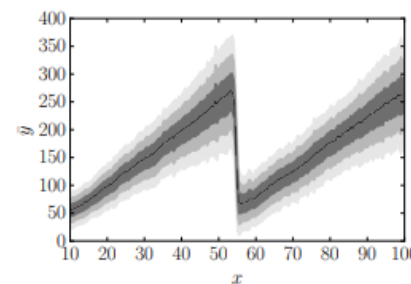
(b) Ground-truth function



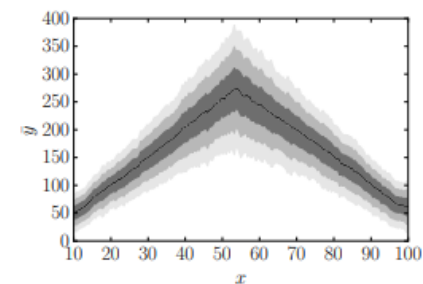
(c) $p_d = 0.2$



(d) $p_d = 0.2$



(e) $p_d = 0.5$



(f) $p_d = 0.5$

불확실도 측면 Trustworthy 구성 모델 대체 AI 모델 개발

□ 불확실도 정보가 해석 가능하며 타당한 AI 프레임워크

- Pruned Diverse Deep Ensemble (Pdiv-DE)
 - Pruning (비슷한 모델 제거) + Diversification (서로 다른 모델)

기존 Deep Ensemble

$$U_{epi}(x) = \frac{1}{M} \sum_{m=1}^M (\mu_m(x) - \bar{\mu}(x))^2 \quad \ast \quad \bar{\mu}(x) = \frac{1}{M} \sum_{m=1}^M \mu_m(x)$$

Submodel들은 유사한 예측함수를 가질 수 있음.

$$\mathcal{E} = \{\mu_1, \mu_2, \dots, \mu_M\}$$

$$\mu_i(x) \approx \mu_j(x)$$

$$\mu_m(x) = w^T h(x) \quad \bar{\mu}(x) = \bar{w}^T h(x)$$

$$\mu_m(x) - \bar{\mu}(x) = (w^T - \bar{w}^T) h(x)$$

비슷한 모델로 data sparsity에 반응하지 않을 수 있음 + 예측 값 크기에 비례

$$U_{epi}^{DE}(x) \sim \mu(x)^2$$

Pruning (예측 값간 차이 유도)

$$\mathcal{E}_p = \{\mu_m \in \mathcal{E} \mid D(\mu_M, \mathcal{E}) > \tau\}$$

Ensemble member간 차이 Pruning threshold

$$U_{epi}^p(x) = \frac{1}{M_p} \sum_{m=1}^{M_p} (\mu_m^p(x) - \bar{\mu}^p(x))^2$$

Diversification (다른 예측 함수 - 활성화함수)

$$\mu_i^p(x) \neq \mu_j^p(x)$$

Data-rich $x \in \mathbf{D}_{dense} \Rightarrow \mu_i(x) \approx \mu_j(x) \Rightarrow U_{epi}(x) \downarrow$

Data-sparse $x \in \mathbf{D}_{sparse} \Rightarrow \mu_i(x) \neq \mu_j(x) \Rightarrow U_{epi}(x) \uparrow$

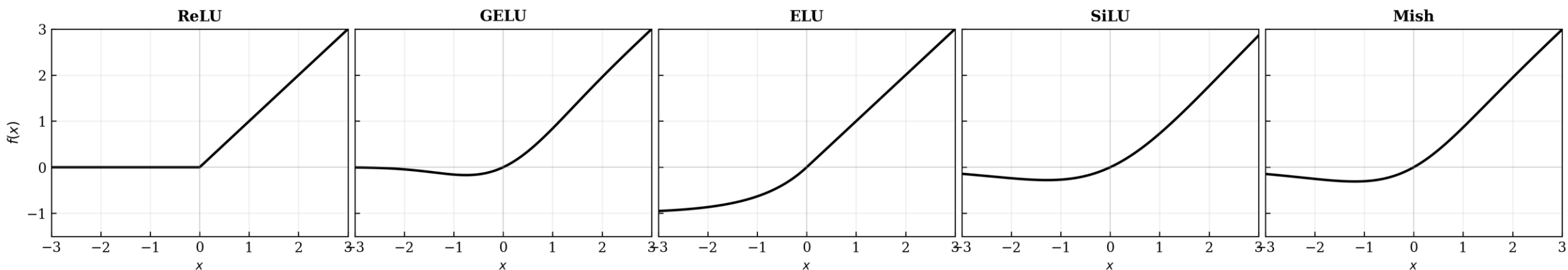
$$U_{epi}^{PDiv-DE}(x) \sim \text{Disagreement}(\mu_1(x), \dots, \mu_M(x))$$

불확실도 측면 Trustworthy 구성 모델 대체 AI 모델 개발

□ 불확실도 정보가 해석 가능하며 타당한 AI 프레임워크

- Pruned Diverse Deep Ensemble (PDiv-DE)
 - Pruning (비슷한 모델 제거) + Diversification (서로 다른 모델)
- 회귀 feature가 일반적 DE보다 더욱 다양함
- 불확실도의 예측 결과에 대한 의존성 완화

Activation	Extrapolation Behavior	Boundary Characteristic
ReLU	Linear (unbounded)	Sharp transition at zero
GELU	Smooth approximation of ReLU	Gradual transition
ELU	Saturates to $-\alpha$ for negative inputs	Bounded below
SiLU (Swish)	Non-monotonic near zero	Smooth, slight dip
Mish	Similar to SiLU, smoother	Self-regularizing



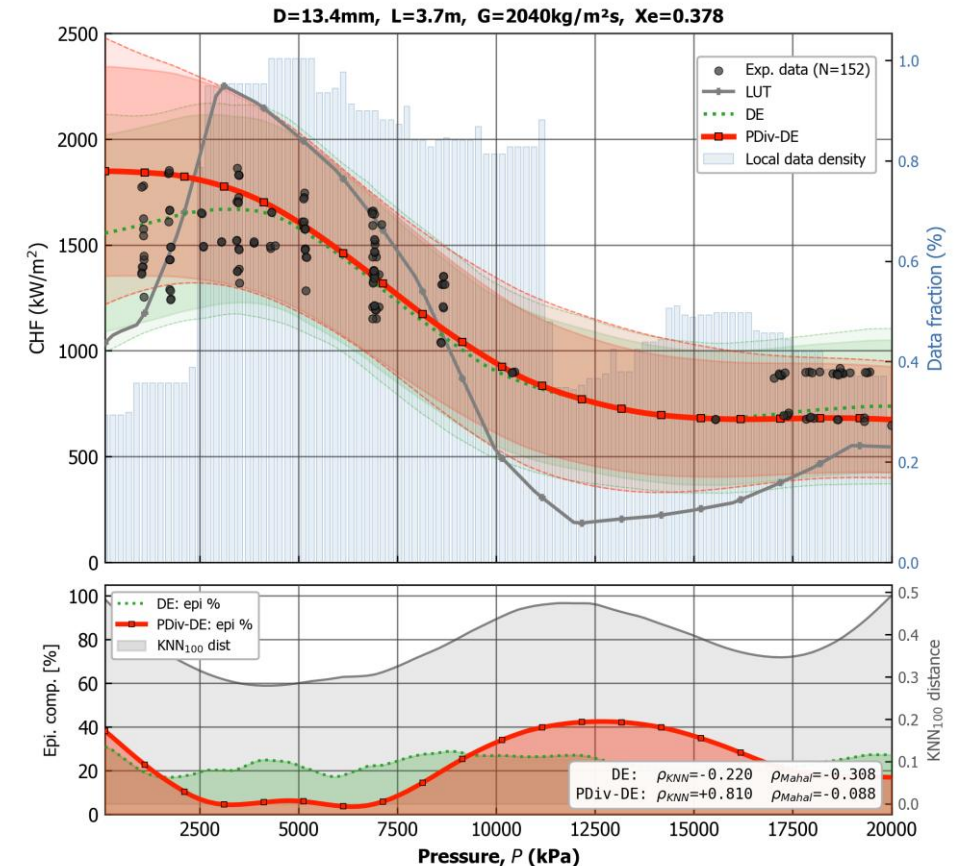
불확실도 측면 Trustworthy 구성 모델 대체 AI 모델 개발

□ 불확실도 정보가 해석 가능하며 타당한 AI 프레임워크

- Pruned Diverse Deep Ensemble (PDiv-DE)
 - 회귀 성능은 Standalone 모델 대비 떨어지지만 기존 예측 방법론보다 우수
 - **Data sparsity와 Epistemic uncertainty간 비례 (정합한 결과)**

Indices	LUT (DSM)	BNN	MCD	DE	PDiv-DE
Mean P/M [-]	1.032	1.102	1.080	1.050	1.088
Std. P/M [-]	0.362	0.341	0.275	0.246	0.325
MAPE [%]	19.77	17.82	15.11	13.24	17.63
RMSPE [%]	36.30	35.59	28.61	25.07	33.71
R ² [-]	0.941	0.961	0.969	0.974	0.959

Indices	Mean ρ_{KNN} for each slice	$N_{slice}(\rho_{KNN} > 0) / N_{total}$	Mean ρ_{Mahal} for each slice	$N_{slice}(\rho_{Mahal} > 0) / N_{total}$
BNN	+0.076	5/10	+0.142	5/10
MCD	+0.417	7/10	+0.540	9/10
DE	+0.328	8/10	0.368	8/10
PDiv-DE	+0.669	10/10	+0.592	9/10

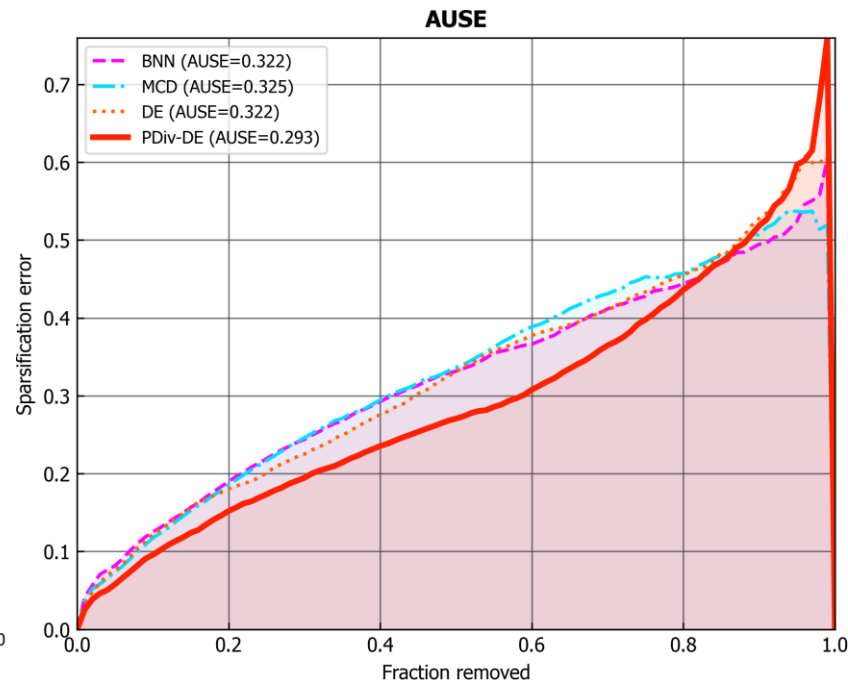
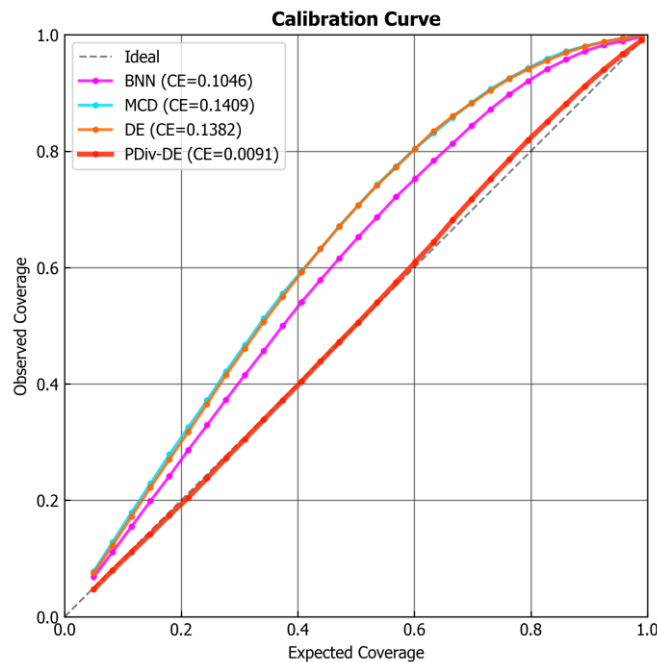


불확실도 측면 Trustworthy 구성 모델 대체 AI 모델 개발

□ 불확실도 정보가 해석 가능하며 타당한 AI 프레임워크

- Pruned Diverse Deep Ensemble (PDiv-DE)
 - ECE가 이상적인 결과에 근접함 (예측 uncertainty가 실제 error에 근접)
 - 낮은 AUSE → Uncertainty에 기여하는 인자를 적절히 반영하고 있음

Indices	BNN	MCD	DE	PDiv-DE
CE	0.104	0.141	0.138	0.009
AUSE	0.322	0.325	0.322	0.293



결론 및 제언

- 인공지능의 정확성, 효율성, 포착 성능에 기반하여 원자로 안전해석에 대한 적용 연구가 증가하고 있음.
- 그러나, 실제 활용을 위해서는 AI의 불확실도에 대해 정량적으로 평가 및 설명할 필요가 있음.
 - 기존 결정론적 모델에서 벗어나 예측과 “불확실도” 정보를 제공가능한 확률론적 모델의 활용이 필요.
 - 확률론적 모델의 불확실도 정보에 대한 신뢰도는 AUSE, ECE (내삽영역), KNN 거리-Mahal 거리와 불확실도간 상관계수(외삽영역)로 평가 가능함.
 - 일반적 확률론적 모델이 제공하는 불확실도 정보는 데이터의 sparsity 및 데이터 자체 분산과의 관계성을 직접적으로 설명하기 어려움.
- 실제 안전해석에 대한 AI 적용을 위해 불확실도와 데이터간 상관관계를 강건하게 설명할 수 있는 확률론적 모델(ex. PDiv-DE) 개발이 필요함.