# LWRS Program Research on Risk Assessment of Safety-related DI&C Systems
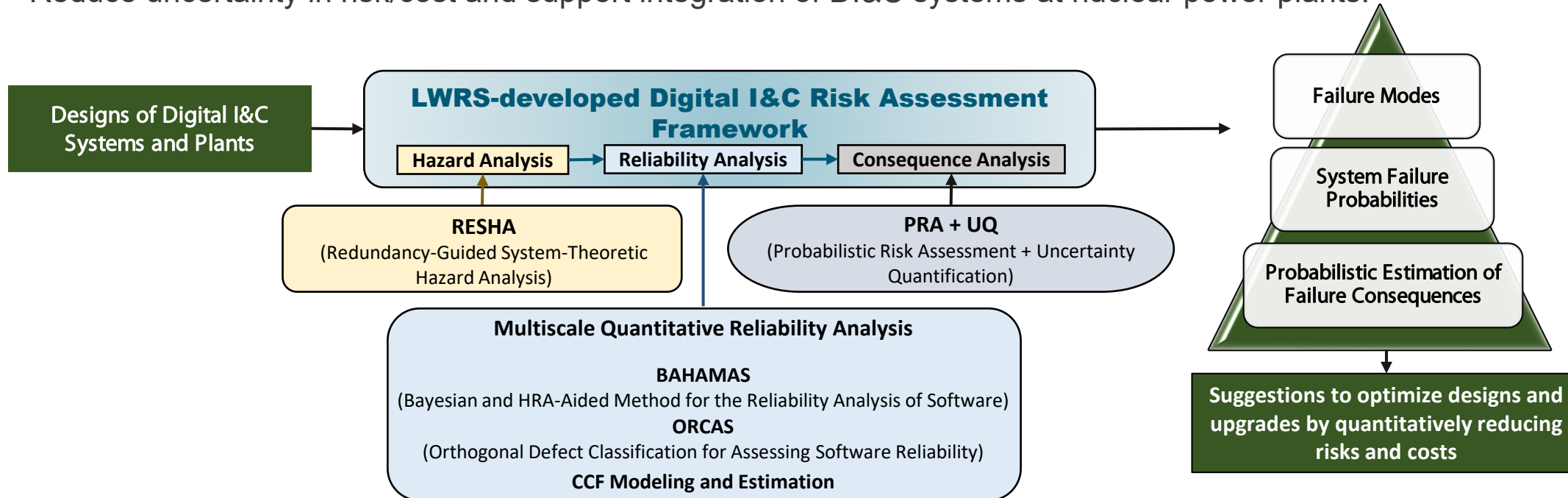
*Project "Digital I&C Risk Assessment"*

*U.S. DOE Light Water Reactor Sustainability (LWRS) Program, Risk-Informed Systems Analysis (RISA) Pathway*

**Congjian Wang, Tate Shorthill, Edward Chen, Jisuk Kim, Jooyoung Park, Svetlana Lawrence**
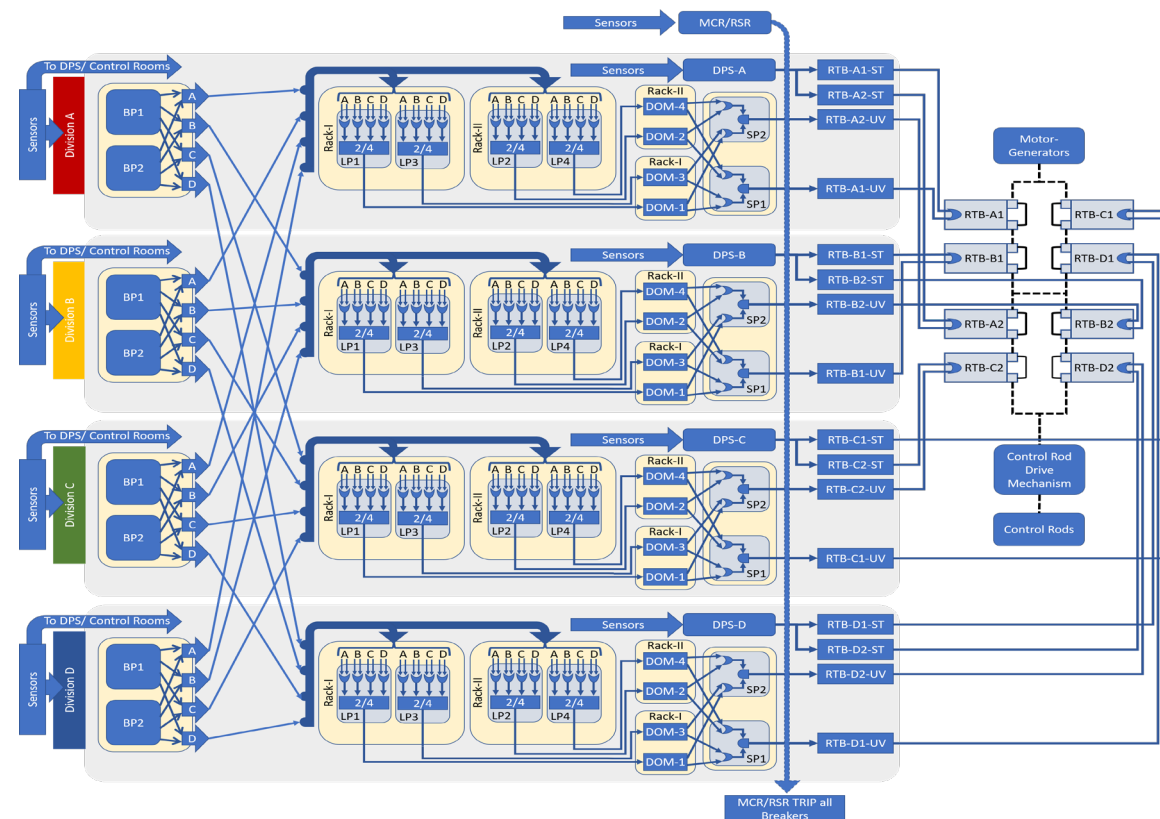
**INL/MIS-24-81813**

November 13, 2024

# Overview of Digital I&C Risk Assessment within LWRS RISA Pathway

- Offer a capability of design architecture evaluation of various DI&C systems to support system design decisions on diversity and redundancy applications

- Develop systematic and risk-informed tools to address CCFs and quantify corresponding failure probabilities for DI&C technologies

- Support and supplement existing risk-informed DI&C design guides by providing quantitative risk-informed and performance-based evidence

- Reduce uncertainty in risk/cost and support integration of DI&C systems at nuclear power plants.

Designs of Digital I&C Systems and Plants

**LWRS-developed Digital I&C Risk Assessment Framework**

Hazard Analysis → Reliability Analysis → Consequence Analysis

**RESHA**
(Redundancy-Guided System-Theoretic Hazard Analysis)

**PRA + UQ**
(Probabilistic Risk Assessment + Uncertainty Quantification)

**Multiscale Quantitative Reliability Analysis**

**BAHAMAS**
(Bayesian and HRA-Aided Method for the Reliability Analysis of Software)
**ORCAS**
(Orthogonal Defect Classification for Assessing Software Reliability)
**CCF Modeling and Estimation**

Failure Modes

System Failure Probabilities

Probabilistic Estimation of Failure Consequences

**Suggestions to optimize designs and upgrades by quantitatively reducing risks and costs**
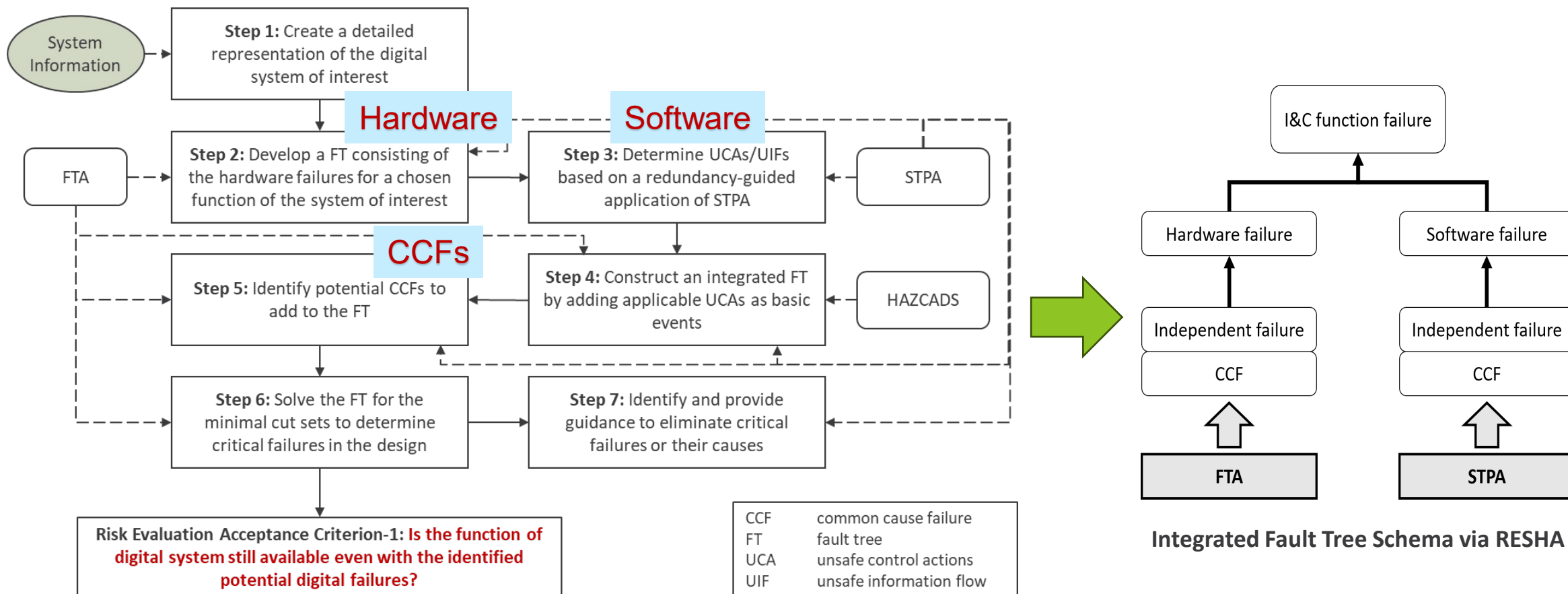
# Value Proposition

- **The framework** is envisioned and developed as **an integrated risk-informed tool** to support vendors and utilities with optimization of design solutions from economical perspectives GIVEN the constrain of meeting risk-informed safety requirements.

- **Quantitative Risk Analysis**
  - Software reliability metrics → DI&C system reliability → Plant safety analysis

- **Risk-informed Design**
  - Management strategy of CCFs
    - Identification and elimination
  - Level of redundancy
    - 4 divisions vs. 2 divisions
    - 4 vs. 2 local logic processors per division
  - Level of diversity
    - Design: Analog? Digital? A combination of both?
    - Software: Design requirements, programming language, etc.
    - Hardware Equipment: Manufacturers, designs, architectures, etc.



*A Four-Division Digital Reactor Trip System*

# Hazard Analysis via Redundancy-guided System-theoretic Hazard Analysis (RESHA)

- Incorporates **Fault Tree Analysis (FTA)** and **System-Theoretic Process Analysis (STPA)**.

- Reframes STPA in a **redundancy-guided** way to **identify various CCFs**.

- Identifies and traces failures in both **Unsafe Control Actions** (UCAs) and **Unsafe Information Flows** (UIFs).



**Integrated Fault Tree Schema via RESHA**

# Quantitative Software Reliability Analysis

- **Methods developed within this project:**

  - **BAHAMAS (**Bayesian and HRA-Aided Method for the Reliability Analysis of Software)

    - Developed for the conditions with limited testing/operational data or for reliability estimations of software in early development stage.

    - Provide an estimation of failure probabilities to support the design of software and target DI&C systems.

  - **ORCAS** (Orthogonal Defect Classification for Assessing Software Reliability)

    - Developed for the conditions with sufficient testing/operational data.

    - A more refined estimation of software failure probabilities can be provided.

| | BAHAMAS | ORCAS |
|---|---|---|
| **Applicable conditions** | • Limited testing/operational data<br>• For reliability estimations of software in early development stage | • Sufficient testing/operational data<br>• For reliability estimations of software in development or testing stage |
| **Key assumption** | Software failures can be traced to human errors in the software development life cycle | Sufficient data is available through testing (e.g., T-Way testing) |
| **Ways to identify root causes** | STPA + ODC + HRA in SDLC | STPA + ODC |
| **Ways to quantify failure rates of root causes** | HRA in SDLC, i.e., Technique for Human Error Rate Prediction | Software reliability growth modeling |

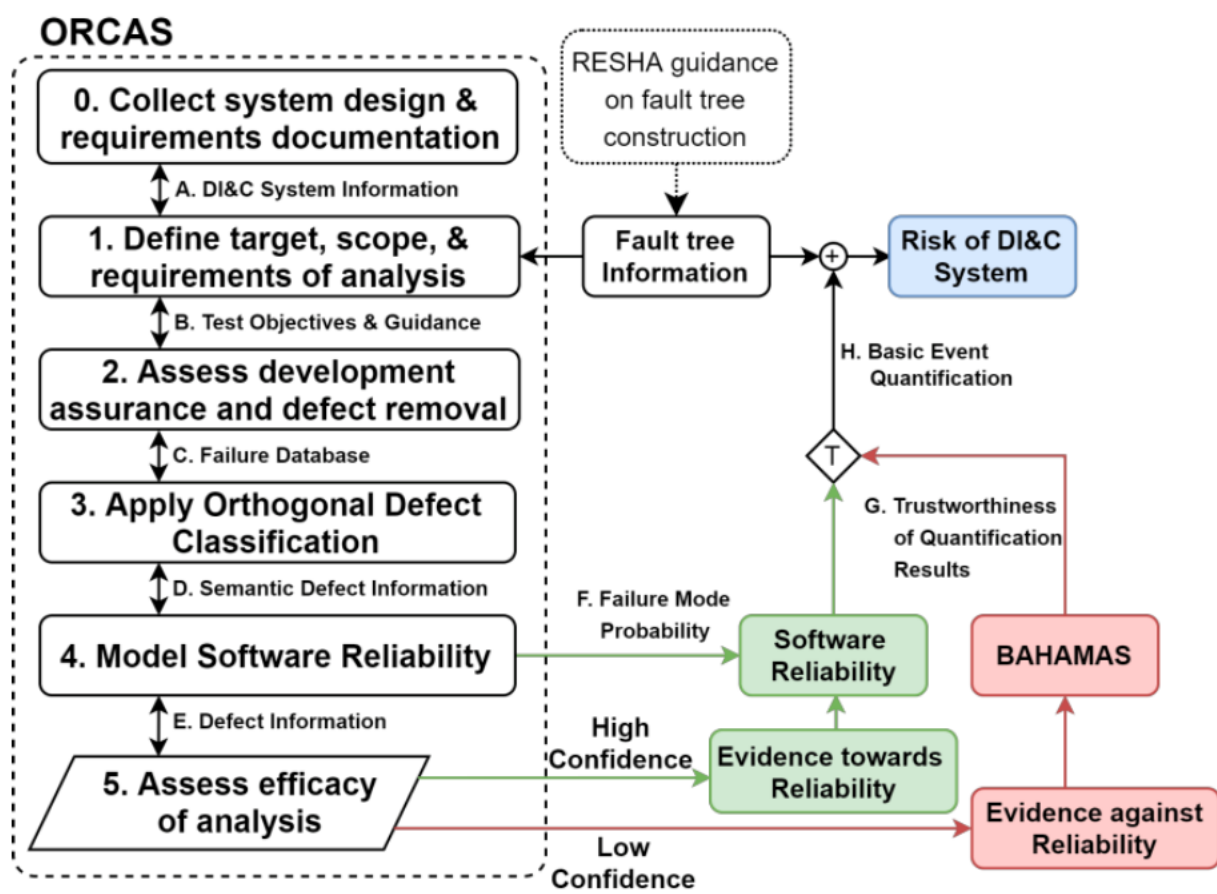| | |
|---|---|
| BNN | Bayesian Belief Network |
| ODC | Orthogonal Defect Classification |
| HRA | Human Reliability Analysis |
| SDLC | software development life cycle |

# Bayesian and HRA-Aided Method for the Reliability Analysis of Software

- **BAHAMAS** tracks human errors in the software development and their influence on the existence of specific types of defects which ultimately influence the probability of software failure
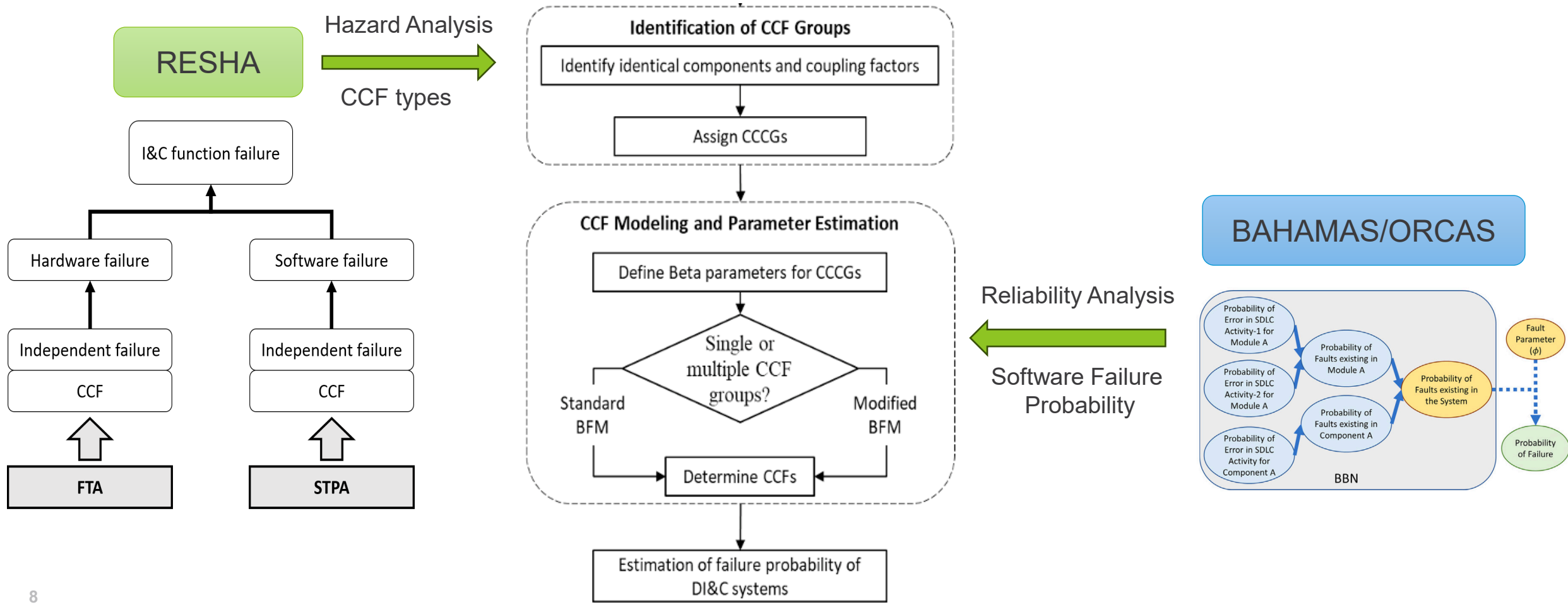
# Orthogonal Defect Classification for Assessing Software Reliability

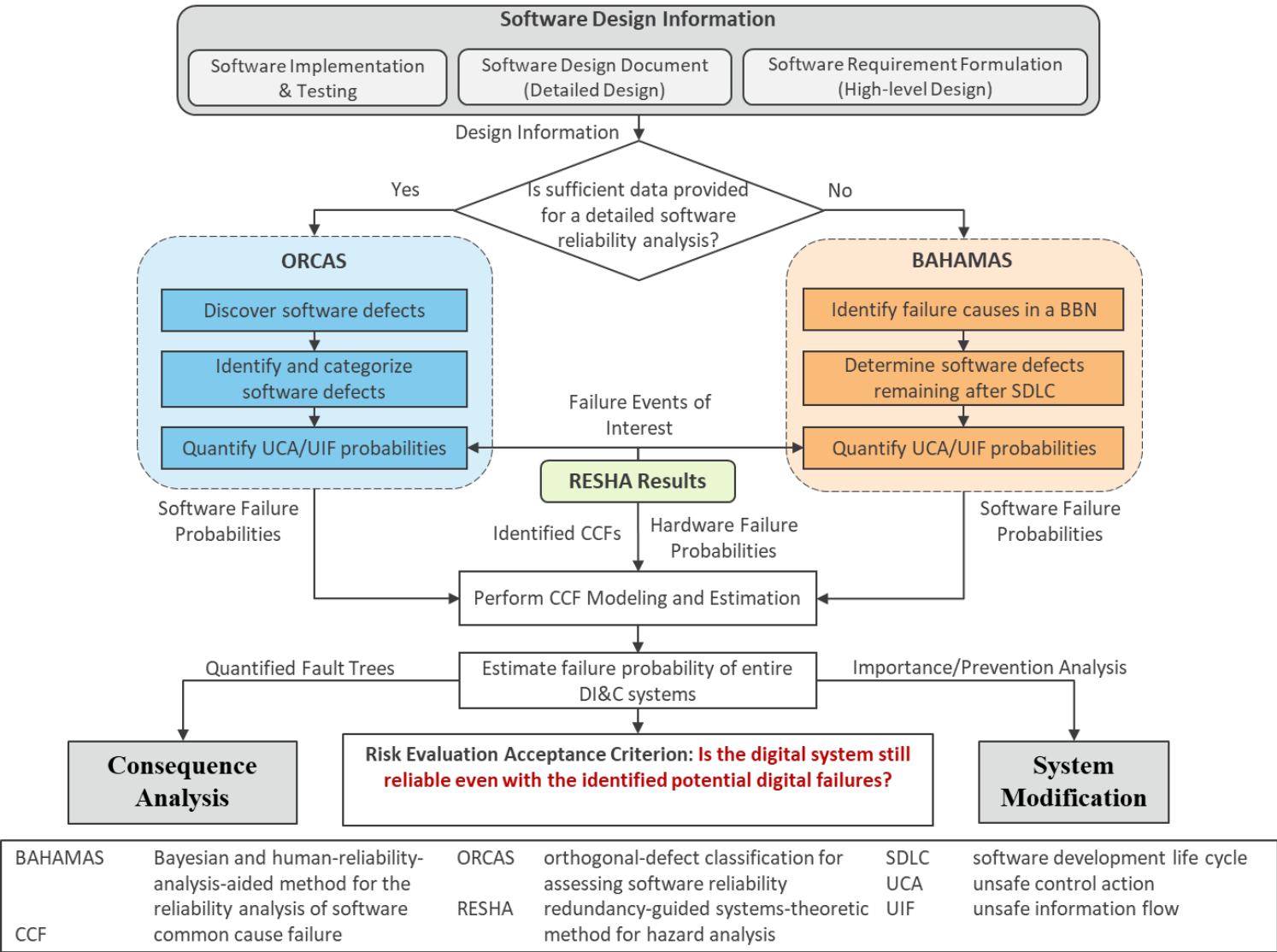- **ORCAS** leverage software comprehensive testing, ODC and software reliability growth models to quantify the software failure probability of specific UCAs/UIFs

# CCF Modeling and Estimation

- A CCF modeling flowgraph is developed for software CCF modeling and estimation based on modified **Beta Factor Model (BFM)** and **Partial Beta Factor (PBF)** Model.

# Multiscale Quantitative Reliability Analysis



**Software Design Information**

| Software Implementation & Testing | Software Design Document (Detailed Design) | Software Requirement Formulation (High-level Design) |

Design Information

Is sufficient data provided for a detailed software reliability analysis?

Yes — **ORCAS**
- Discover software defects
- Identify and categorize software defects
- Quantify UCA/UIF probabilities

No — **BAHAMAS**
- Identify failure causes in a BBN
- Determine software defects remaining after SDLC
- Quantify UCA/UIF probabilities

Failure Events of Interest

**RESHA Results**

Software Failure Probabilities · Identified CCFs · Hardware Failure Probabilities · Software Failure Probabilities

Perform CCF Modeling and Estimation

Quantified Fault Trees — Estimate failure probability of entire DI&C systems — Importance/Prevention Analysis

**Consequence Analysis**

**Risk Evaluation Acceptance Criterion: Is the digital system still reliable even with the identified potential digital failures?**

**System Modification**

| | | | | | |
|---|---|---|---|---|---|
| BAHAMAS | Bayesian and human-reliability-analysis-aided method for the reliability analysis of software | ORCAS | orthogonal-defect classification for assessing software reliability | SDLC | software development life cycle |
| | | | | UCA | unsafe control action |
| CCF | common cause failure | RESHA | redundancy-guided systems-theoretic method for hazard analysis | UIF | unsafe information flow |

9

# Major Accomplishments in FY-24 (I)

➢ Developed a novel approach to evaluate the reliability of ML-integrated control systems.

- A journal article published.
- Results were included in the June M4 technical report.

# Major Accomplishments in FY-24 (II)

➢ Collaboration with PWROG for CCFs quantification for DI&C system
  - Results were included in a proprietary technical white paper.

➢ Initiated collaboration with GE Hitachi for function-based risk assessment of multi-function DI&C systems.
  - Results were included in August M3 technical report.



INL/RPT-24-04807

**Light Water Reactor Sustainability Program**

**Methodology Illustration:**
**Qualitative Risk Analysis of the GEH**
**C10 Safety System**

**OFFICIAL USE ONLY**

May be exempt from public release under the Freedom of Information Act (5 U.S.C. 552), exemption number and category: Exemption No. 4 Commercial Proprietary.
Department of Energy review required before public release.
Name/Org: Edward Chen  Date: 7/30/2024
Guidance (if applicable): N/A

**PROPRIETARY INFORMATION**

This document contains Proprietary Information disclosed under and in accordance with the Non-Disclosure Agreement (NDA) No. 23NDA228 Rev. 0 between Battelle Energy Alliance, LLC (BEA), GE-Hitachi Nuclear Energy Americas LLC (GEH) dated 06/07/2023, and is not to be further disclosed by anyone gaining access to this document as a result of any disclosure under such NDA, without the prior wr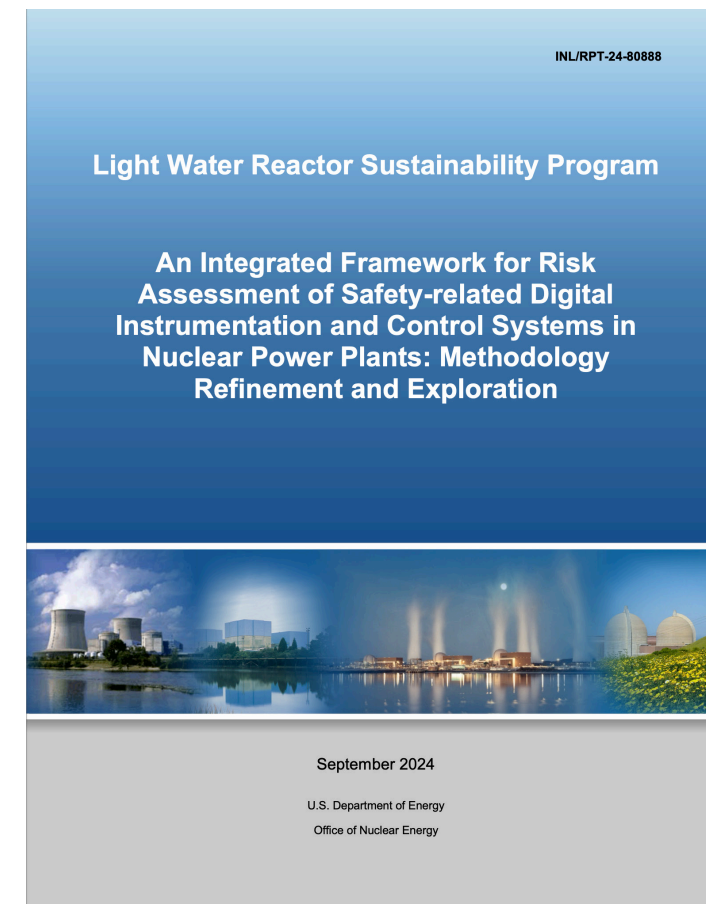itten approval of BEA, and GEH, except as expressly provided for in the above stated NDA. Recipient is responsible for compliance with United States laws and regulations governing export controls including, but not limited to, the Export Administration Regulations (EAR) (15 CFR Parts 730-774), the International Traffic in Arms Regulations (ITAR) (22 CFR Parts 120-130), and the Nuclear Regulatory Commission and Department of Energy export regulations (10 CFR Parts 110 and 810). Unauthorized export, deemed export, or re-export without an export license may result in administrative, civil, or criminal penalties.

August 2024

U.S. Department of Energy
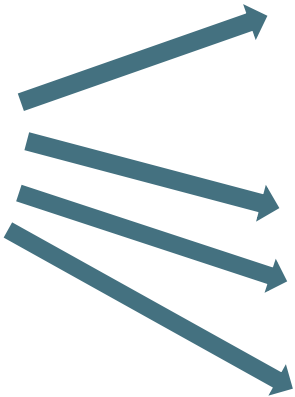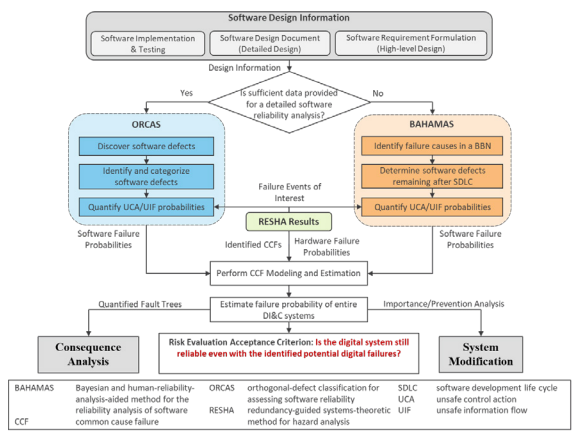Office of Nuclear Energy

**OFFICIAL USE ONLY**

# Major Accomplishments in FY-24 (III)

➢ Refined the reliability analysis methods of a safety-critical DI&C system

- Developed a novel approach to evaluate inter-system CCFs in highly redundant and diverse DI&C systems

- Root cause correlation analysis via ORCAS and natural language processing (NLP)

- Leverage LLM to perform hazard analysis and reliability analysis

- Methodology improvement and a user guidance for industry use was included in the September M2 technical report.

INL/RPT-24-80888

**Light Water Reactor Sustainability Program**

**An Integrated Framework for Risk Assessment of Safety-related Digital Instrumentation and Control Systems in Nuclear Power Plants: Methodology Refinement and Exploration**

September 2024

U.S. Department of Energy

Office of Nuclear Energy

# Roadmap: From Risk Assessment to Design Optimization and Licensing

*PRA: probabilistic risk assessment

# Roadmap: Risk Assessment Framework Development

- Software for Hazard Identification and Evaluation of Digital Systems (SHIELDS)

# Research Activities in FY-25

- Improve and further develop the current methods for risk assessment of multi-function DI&C systems
  - Keep supporting the need of DI&C reliability analysis from the industry.

- Refine the current methods:
  - Intra- and inter-system CCF modeling
  - Align better with international standards and existing risk-informed approaches and guides.

- SHIELDS framework development

- Develop capabilities on risk-informed evidence generation and evaluation to support DI&C safety assurance and design optimization with the industry and other research institutions.

- Develop novel approaches to inform risk management and design optimization of advanced (semi-) autonomous DI&C systems designed for existing LWR fleets.

# Collaborators

- **PWROG engagement**: Digital I&C reliability analysis and CCF evaluation

- **GE Hitachi**: function-based risk assessment of multi-function DI&C platforms

- **KAERI:** safety analysis of advanced DI&C technologies

- **NEI/Halden**: DI&C D3 application and safety assurance

- **NRC**: Risk evaluation and design optimization of AI-aided DI&C systems

- **North Carolina State University**:

    – DI&C hazard analysis using large language model

    – CCF analysis and parameter estimation using model-based approaches.

# Sustaining National Nuclear Assets

*lwrs.inl.gov*