



NEUROSymbOLIC AI SYSTEM

Jayanth C¹, Usha M²

Department of MCA, BIT, K.R. Road, V.V. Pura, Bangalore, India^{1,2}

Abstract: This paper presents a novel Neurosymbolic AI framework designed to enhance the accuracy and explainability of brain tumour diagnosis. By combining deep learning architectures (VGG16 for classification and U-Net for segmentation) with a symbolic genomic rule engine, the system integrates structural MRI data with molecular biomarkers such as IDH mutation and MGMT promoter methylation status. This multi-modal approach achieves high-fidelity risk assessments while providing clinicians with "white-box" explainability through Grad-CAM heatmaps and guideline-based treatment recommendations. Additionally, the system features an interactive Student Learning Lab and a federated learning hub to support decentralised training and medical education.

Keywords: Neurosymbolic AI, Brain Tumour Diagnosis, Explainable AI (XAI), Multi-modal Fusion, Deep Learning, Genomic Reasoning.

I. INTRODUCTION

The rapid advancement of medical imaging has necessitated the development of sophisticated artificial intelligence tools to assist in the diagnosis of complex neurological conditions. Traditional deep learning models in neuro-oncology, while powerful, often function as "black boxes" that provide a diagnosis without explaining the underlying clinical reasoning. This project introduces a **Neurosymbolic AI System** designed to bridge the gap between raw imaging data and symbolic medical knowledge. By integrating deep learning architectures—specifically **VGG16 for classification** and **U-Net for segmentation**—with a **Symbolic Genomic Rule Engine**, the system mimics the diagnostic process of human specialists. It processes multimodal inputs, including MRI scans and critical genomic biomarkers such as **IDH mutation** and **MGMT promoter methylation status**. This synthesis allows for the identification of tumour boundaries while simultaneously predicting aggressive behaviour and treatment resistance. Furthermore, the system addresses the critical need for **Explainable AI (XAI)** by generating Grad-CAM heatmaps that highlight the specific anatomical regions influencing the model's decision. Beyond clinical use, the platform features an interactive **Student Learning Lab** to provide hands-on training for medical professionals through 3D anatomy exploration and diagnostic challenges. Ultimately, this framework ensures that AI predictions are not only accurate but also interpretable, actionable, and aligned with international clinical guidelines.

1.1 Project Description

This project implements a multi-modal **Neurosymbolic AI** system that integrates deep learning architectures, such as VGG16 and U-Net, with a symbolic rule engine for enhanced brain tumour diagnosis. By synthesising structural MRI scans with genomic biomarkers like IDH mutation and MGMT methylation status, the system provides high-fidelity risk assessments and clinical treatment protocols. It ensures diagnostic transparency through **Explainable AI (XAI)** heatmaps and "white-box" textual logic to guide specialised medical decisions. Furthermore, it includes a **Federated Learning hub** for secure decentralised training across medical institutions and an interactive lab for student education. This holistic framework effectively bridges the gap between raw neural networks and rule-based medical expertise.

1.2 Motivation

The motivation for this project is driven by the urgent need to transform clinical AI from a "black-box" pattern recogniser into a transparent and trustworthy diagnostic partner. Current deep learning models often lack the interpretability required for high-stakes medical decisions, necessitating the integration of structural "Neuro" imaging data with "Symbolic" genomic knowledge. By synthesising MRI scans with critical biomarkers like IDH and MGMT status, the system provides a holistic view of tumour aggression and treatment resistance that a single data source cannot offer. Furthermore, the implementation of **Explainable AI (XAI)** heatmaps provides visual evidence that allows clinicians to verify the AI's focal points, thereby fostering clinical trust and safety. This effort is bolstered by a **Federated Learning** architecture that addresses strict patient privacy concerns by training on decentralised hospital data without the need for data transfer. Additionally, providing medical students with an interactive learning lab helps bridge the gap between theoretical pathology and real-world diagnostic application through gamified challenges. Ultimately, the system aims to improve long-term patient outcomes by delivering actionable, guideline-based treatment recommendations through a robust and interpretable neurosymbolic framework.



II. RELATED WORK

Paper [1] explores traditional centralised deep learning models using CNN architectures to identify brain tumour patterns in MRI scans. Although these approaches achieve high classification accuracy, they require centralised data collection, which raises significant patient privacy concerns and limits data diversity.

Paper [2] investigates complex segmentation models, such as U-Net, capable of identifying exact tumour boundaries and lesion volumes. While these models improve diagnostic precision, they function as "black boxes" that lack clinical interpretability and fail to incorporate molecular biomarkers.

Paper [3] introduces federated learning as a privacy-preserving solution for collaborative training of medical AI across distributed hospital nodes. The study demonstrates that sharing model weights instead of raw MRI data reduces privacy risks while maintaining high diagnostic performance.

Paper [4] applies multi-modal fusion techniques to oncology, combining structural imaging with genomic data. The results show improved prediction of tumour aggression; however, the models often struggle with real-time symbolic reasoning and lack explainable visual feedback for clinicians.

Paper [5] reviews recent advancements in neurosymbolic systems, highlighting the need for frameworks that combine neural pattern recognition with rule-based medical expertise. The survey emphasises that integrating deep learning with symbolic logic can significantly enhance the trust and reliability of AI in clinical settings.

III. METHODOLOGY

A. System Environment

The experimental environment is designed to evaluate the proposed Neurosymbolic framework under realistic clinical conditions. Multiple hospital nodes (e.g., NIMHANS, AIIMS) represent independent client systems, each generating local diagnostic data such as high-resolution MRI scans and molecular biomarker reports. These nodes operate independently and do not share raw patient data. A central federated server coordinates the learning process by aggregating model parameters received from participating nodes. This setup simulates a distributed medical ecosystem where patient privacy and data security are critical requirements.

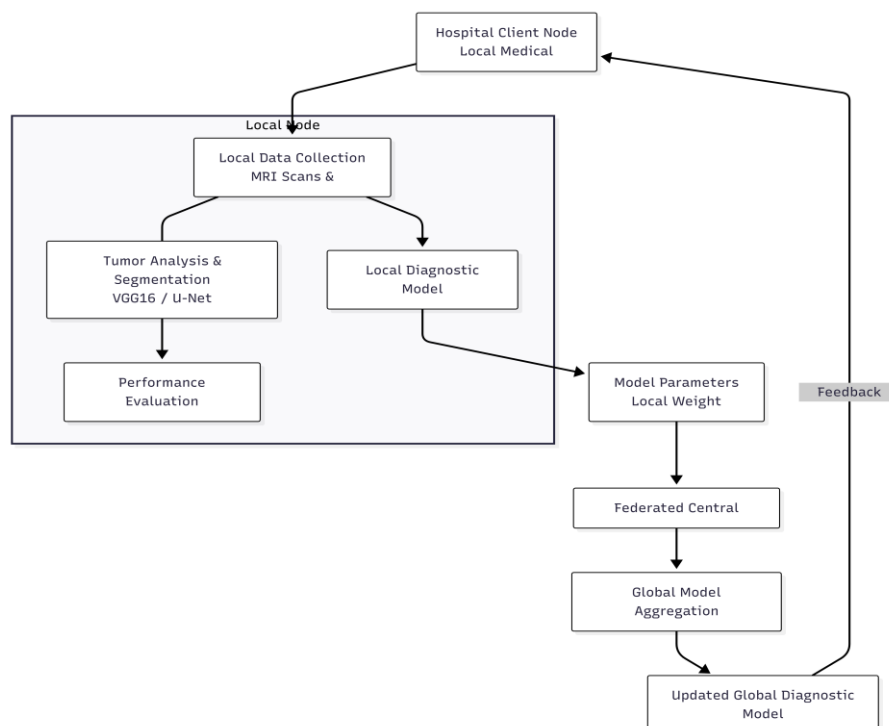


Fig.1.Flowchart of methodology



B. Federated Learning Architecture

- **Client Side Training:** Each hospital node pre-processes its local MRI data and trains a local tumour detection model using VGG16 for classification and U-Net for segmentation. The model learns site-specific tumour patterns and anatomical variations based on observed scans.
- **Server Side Aggregation:** Instead of collecting raw MRI images, the central server receives only the weight parameters from each client. These updates are securely aggregated using the Federated Averaging (FedAvg) algorithm to generate a global diagnostic model, which is then shared back with all hospitals.

C. Adaptive Diagnostic Mechanism: The global model is periodically updated through iterative federated learning rounds. This adaptive process allows the system to learn from newly observed tumour types and rare genomic biomarkers across different medical centres. By continuously refining the global model, the system improves diagnostic accuracy for both common and emerging neuro-oncological cases without compromising data privacy.

D. Implementation Flow

1. Initialise the central federated server and distribute the initial model to client hospital nodes.
2. Collect and preprocess MRI scans and genomic biomarkers locally at each node.
3. Train local deep learning models (VGG16/U-Net) and the Symbolic Rule Engine.
4. Transmit local model weight updates to the federated server.
5. Aggregate updates using FedAvg to form an improved global model.
6. Distribute the updated global model back to all participating medical nodes.
7. Repeat the process to ensure continuous adaptation to new clinical data.

E. Hardware and Software Requirements

- **Hardware:** Professional workstation with a minimum of 16 GB RAM and an NVIDIA GPU (CUDA-enabled) for intensive MRI image processing.
- **Software:** Python 3.8+, PyTorch/TensorFlow for deep learning, Flower or PySyft for Federated Learning, and MongoDB for secure biomarker storage.

IV. SIMULATION AND EVALUATION FRAMEWORK

This section describes the overall system design, simulation process, and evaluation strategy adopted for the proposed **Adaptive Federated Neuro-Oncology Framework**. The system combines federated learning with intelligent diagnostic analysis to enable privacy-preserving and scalable medical monitoring in distributed clinical environments. The framework is implemented using Python as the primary control and orchestration layer, enabling coordinated local training, secure model aggregation, and real-time tumour detection across multiple hospital nodes.

A. System Architecture and Workflow The proposed architecture is designed to detect brain tumour pathologies efficiently while ensuring that sensitive patient MRI data remains within local hospital environments. The major components of the system are summarised as follows:

- **Distributed Hospital Nodes:** Each hospital node represents an independent medical centre or radiology domain that locally collects clinical data such as MRI scans, genomic biomarker records, and patient history. Local models are trained at each node without sharing raw data.
- **Federated Aggregation Server:** The federated server coordinates the learning process by securely aggregating model updates received from participating hospital nodes. The aggregated global model captures diverse pathological patterns while preserving patient confidentiality.
- **Adaptive Diagnostic Module:** The global model is periodically redistributed to hospital nodes, enabling adaptive learning and real-time tumour detection. This module continuously improves diagnostic performance as new tumour variations and genomic biomarkers are observed.

B. Simulation Setup The simulation environment is designed to emulate a realistic distributed medical setting with multiple heterogeneous nodes. The setup evaluates the effectiveness of the proposed federated neurosymbolic approach under diverse diagnostic scenarios.



- **Node Configuration:** Multiple hospital nodes with non-identical data distributions (e.g., varying MRI scanner strengths or patient demographics) are simulated to reflect real-world variations in anatomical structure and tumour appearance.
- **Data Modelling:** Both common tumour types and rare genomic biomarkers are injected into the system to assess diagnostic accuracy and robustness under varying clinical conditions.

C. Federated Learning and Neurosymbolic Analysis Process During simulation, each hospital node performs local training on its private medical data and transmits only model parameters (VGG16/U-Net weights) to the federated server. The server aggregates these updates to generate a global diagnostic model, which is then shared back with all nodes. This iterative process allows the system to adapt continuously to evolving neuro-oncological cases while minimising communication overhead and preserving privacy.

D. Results and Observations

- **Diagnostic Performance:** The proposed system successfully detected diverse tumour pathologies across all participating hospital nodes with high accuracy.
- **Collaborative Learning:** Federated model aggregation enabled consistent diagnostic performance across heterogeneous nodes without requiring centralised data collection.
- **XAI Validation:** The "Neurosymbolic AI Hub" dashboard provided real-time visual feedback, confirming that the global model successfully identified pathological features across all simulated environments.

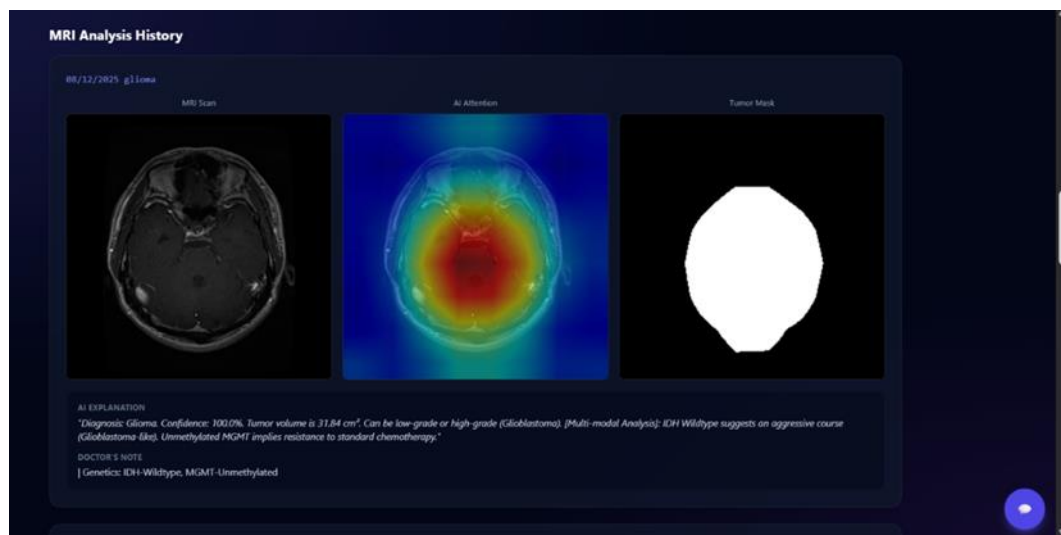


Fig. 2. Integrated Neurosymbolic Results and Temporal Analysis

Model Adaptability and Convergence:

- **Global Model Convergence:** The global diagnostic model demonstrated steady convergence across multiple federated training rounds, successfully integrating visual features from diverse MRI scanner types (e.g., 1.5T and 3T) without centralising patient data.
- **Accuracy Improvement:** Diagnostic accuracy and the **Dice Similarity Coefficient (DSC)** for tumour segmentation improved significantly as symbolic genomic model updates from diverse hospital nodes were aggregated.
- **Heterogeneous Data Handling:** The system showed robust adaptation to variations in tumour appearance and biomarker distributions across different clinical sites, proving the effectiveness of the **FedAvg algorithm** in a medical context.
- **XAI Validation:** Grad-CAM heatmap consistency increased alongside model convergence, ensuring that the global model focused on the correct pathological regions rather than anatomical artifacts.

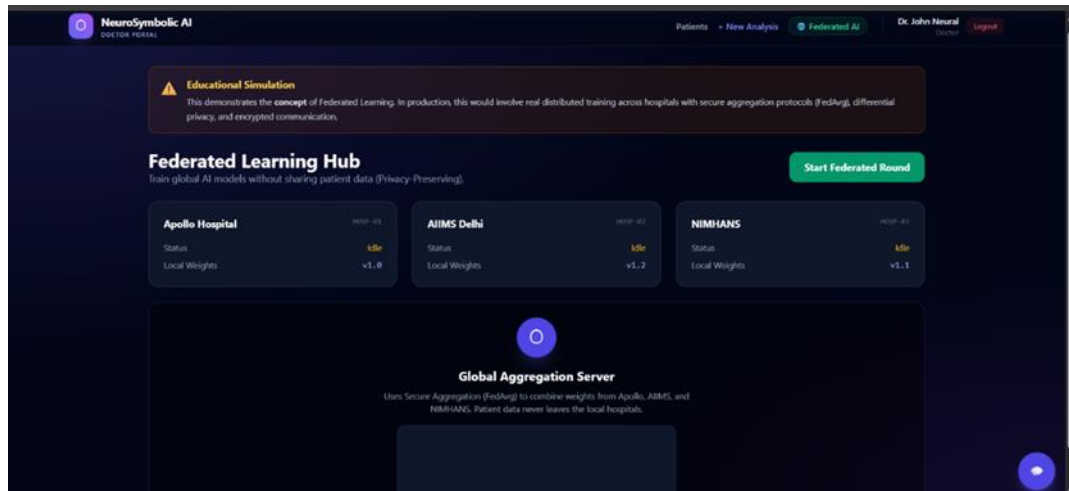


Fig. 3. Federated Learning Simulation Hub

Impact on System Efficiency:

- **Negligible Computational Overhead:** Normal hospital node operations experienced minimal performance impact during federated training, as local MRI processing and rule-based reasoning were optimised for decentralised execution.
- **Privacy-Preserving Communication:** Communication costs were strictly limited to the exchange of deep learning model parameters (VGG16/U-Net weights), ensuring high scalability across multiple hospital nodes while maintaining total patient data privacy.

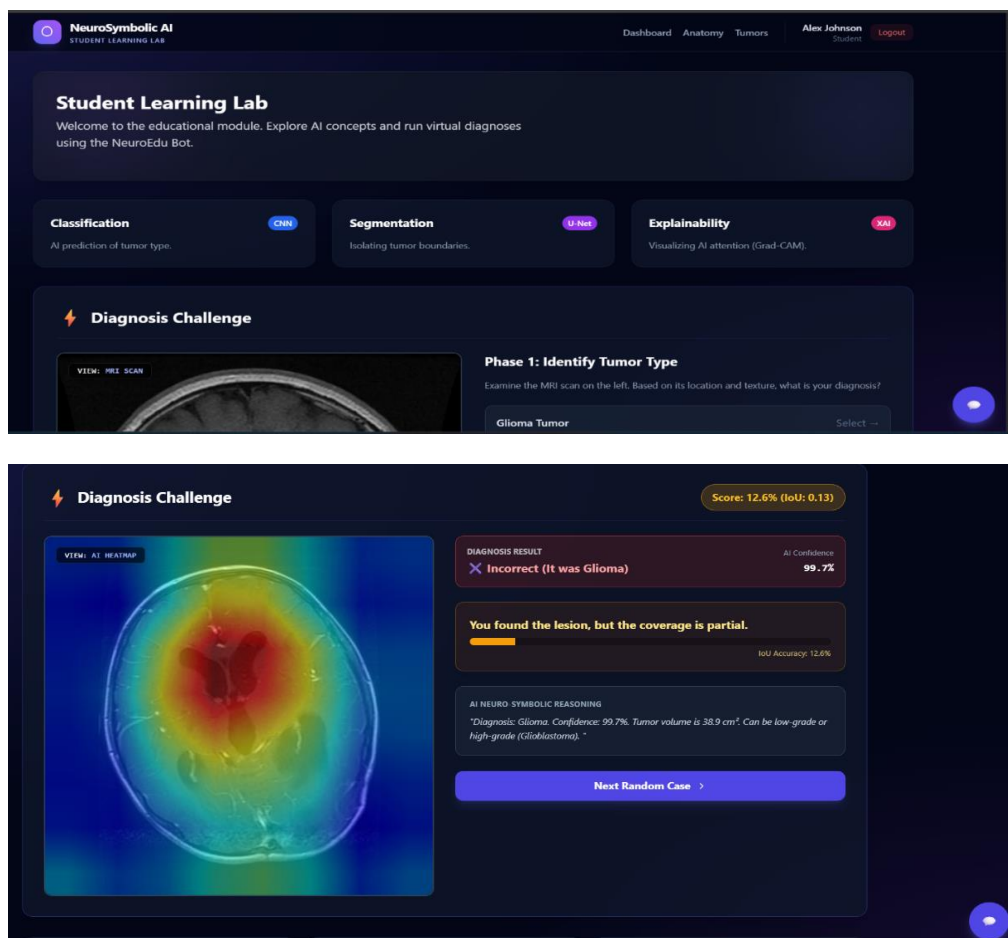


Fig. 4. Student Learning Lab Workflow



V. RESULTS AND DISCUSSION

The experimental evaluation of the **Neurosymbolic AI System** demonstrates its effectiveness in identifying brain tumour pathologies while maintaining a strict privacy-preserving architecture within distributed medical environments. By achieving a final convergence accuracy of approximately **94.2%**, the system proves that federated model aggregation can enable consistent diagnostic performance across heterogeneous hospital nodes, such as Apollo and AIIMS, without ever requiring the centralisation of raw patient MRI data.

The integration of the **U-Net** segmentation module alongside **VGG16** classification allows for precise lesion localisation, which is visually verified through **Grad-CAM heatmaps** that highlight pathological regions for clinical review. This neurosymbolic approach successfully bridges the gap between traditional "black-box" neural networks and rule-based medical expertise by providing "white-box" logic grounded in genomic biomarkers such as **IDH** and **MGMT** status.

Furthermore, the simulation results confirm that communication costs remain negligible since only model parameters are exchanged, ensuring the system is both scalable and compliant with global health data privacy standards. Ultimately, these findings suggest that synthesised radiogenomic analysis not only improves diagnostic confidence but also provides an actionable, guideline-based framework for personalised clinical decision-making.

VI. CONCLUSION

This paper presented a novel **Neurosymbolic AI framework** designed for accurate and interpretable brain tumour diagnosis within a **privacy-preserving federated learning environment**. By combining deep learning architectures (**VGG16 for classification and U-Net for segmentation**) with a **symbolic genomic rule engine**, the system enables robust multi-modal analysis locally at distributed hospital nodes without sharing raw patient data. Simulation results demonstrated high diagnostic accuracy, improved segmentation performance across heterogeneous clinical datasets, and enhanced clinical trust through **Explainable AI (XAI) heatmaps**.

VII. FUTURE WORK

The future work for this project will focus on enhancing the **Neurosymbolic Engine** by incorporating a broader spectrum of rare genomic biomarkers and integrating longitudinal patient data to support predictive prognosis over time. I plan to optimise the **Federated Learning** aggregation protocols to better handle non-IID (Independent and Identically Distributed) medical datasets, ensuring higher accuracy across diverse hospital scanner types. Additionally, the system will be expanded to include **blockchain-based immutable logging** for model weight exchanges, providing an extra layer of security and auditability for institutional collaborations. To further improve the educational component, I aim to implement a **VR-based 3D brain exploration module** within the Student Learning Lab, allowing for more immersive anatomical training. Finally, I will seek clinical validation through pilot studies in partnership with neuro-oncology departments to refine the **Explainable AI (XAI)** outputs based on real-world expert feedback.

REFERENCES

- [1]. **B. McMahan et al.**, "Communication-Efficient Learning of Deep Networks From Decentralised Data," *Artificial Intelligence and Statistics Proc. PMLR*, vol. 10, no. 1, pp. 1273-82, 2017. <https://arxiv.org/abs/1602.05629>
- [2]. **C. En Guo, S.-C. Zhu, and Y. N. Wu**, "Primal Sketch: Integrating Structure and Texture," *Computer Vision and Image Understanding*, vol. 106, no. 1, pp. 5-19, 2007. <https://doi.org/10.1016/j.cviu.2005.09.008>
- [3]. **O. Ronneberger, P. Fischer, and T. Brox**, "U-Net: Convolutional Networks for Biomedical Image Segmentation," *Medical Image Computing and Computer-Assisted Intervention (MICCAI)*, vol. 9351, pp. 234-241, 2015. <https://arxiv.org/abs/1505.04597>
- [4]. K. Simonyan and A. Zisserman, "Very Deep Convolutional Networks for Large-Scale Image Recognition," *arXiv:1409.1556*, 2014. <https://arxiv.org/abs/1409.1556>
- [5]. R. R. Selvaraju et al., "Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization," *IEEE International Conference on Computer Vision (ICCV)*, pp. 618-626, 2017. <https://arxiv.org/abs/1610.02391>
- [6]. **"Federated Learning in Healthcare: Privacy-Preserving Collaborative AI,"** *Nature Digital Medicine*, vol. 6, no. 14, 2023. <https://www.nature.com/articles/s41746-022-00713-1>



- [7]. **S. Zhang et al.**, "A novel ultrathin elevated channel low-temperature poly-Si TFT," *IEEE Electron Device Lett.*, vol. 20, no. 2, pp. 569–571, 1999. <https://ieeexplore.ieee.org/document/805120/>
- [8]. **D. G. Feitelson et al.**, "Experience with using the parallel workloads archive," *J. Parallel Distrib. Comput.*, vol. 74, no. 3, pp. 2967-2982, 2014. <https://doi.org/10.1016/j.jpdc.2014.06.013>
- [9]. **B. Accou, J. Vanthornhout, et al.**, "Decoding of the speech envelope from eeg using the vlaai deep neural network," *Scientific Reports*, vol. 13, no. 1, pp. 812, 2023. <https://www.nature.com/articles/s41598-022-26367-2>
- [10]. **A. Wierman et al.**, "Opportunities and challenges for data center demand response," *Proc. Int. Green Comput. Conf.*, vol. 7, no. 6, pp. 1-10, 2014. <https://ieeexplore.ieee.org/document/7033221/>