

# Standardized ML EDA Data Pipelines with Ontology-driven Data Mapping

Vidya A. Chhabria, Arizona State University

# Outline

- **Challenges and motivation for ML EDA data pipelines**
- Si2 open standards working group ML EDA data pipeline
  - Proposed flow overview
  - Components of the flow with an example usecase
- Conclusion

# Challenges with Data Formats

Challenges with data formats:

- AI-unfriendly data structures and interfaces of EDA tools, hindering high volumes of data queries, transfer, and processing
- Steep learning curve of EDA knowledge and tools for AI experts
- Lack of shared intermediate representation, hindering data reuse across projects

Challenges with creating a dataflow pipelines

- Lack of an ontology to standardize naming convention
- Diverse tools with inconsistent data formats
- EDA tool servers are different from ML GPU servers and require the transfer of large amount of data quickly



# CircuitOps: Existing Dataformat for ML EDA

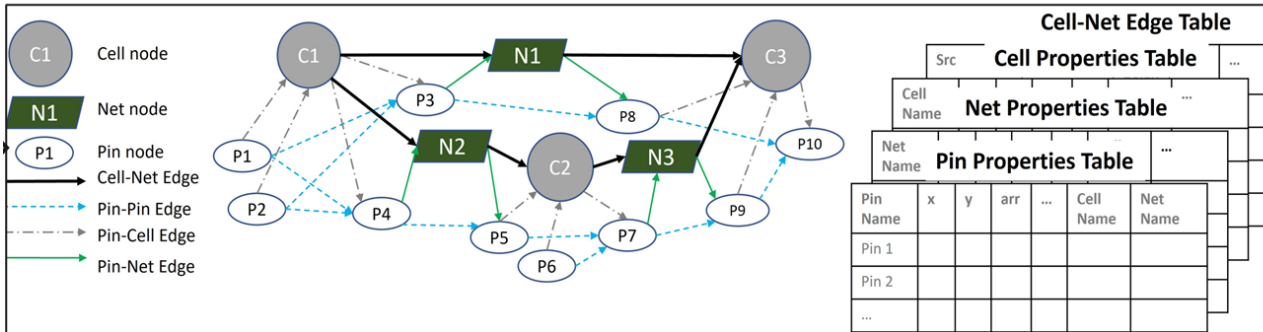
Slide credits: Rongjian Liang, NVIDIA

<https://github.com/NVlabs/CircuitOps>

- Labeled-property graph-based data format backed by IR tables (csvs)
- Suitable for training data generation for different ML EDA applications by filtering out data from LPG

## CircuitOps: ML-friendly data representation format within OpenROAD

pandas.DataFrame features



- Data communication between servers and lack of ontology still remains a challenge

# Outline

- Challenges and motivation for ML EDA data pipelines
- Si2 open standards working group ML EDA data pipeline
  - **Proposed flow overview**
  - Components of the flow with an example usecase
- Conclusion

# Goal for Si2 Open Standards Working Group

## Si2 Open Standards Working Group

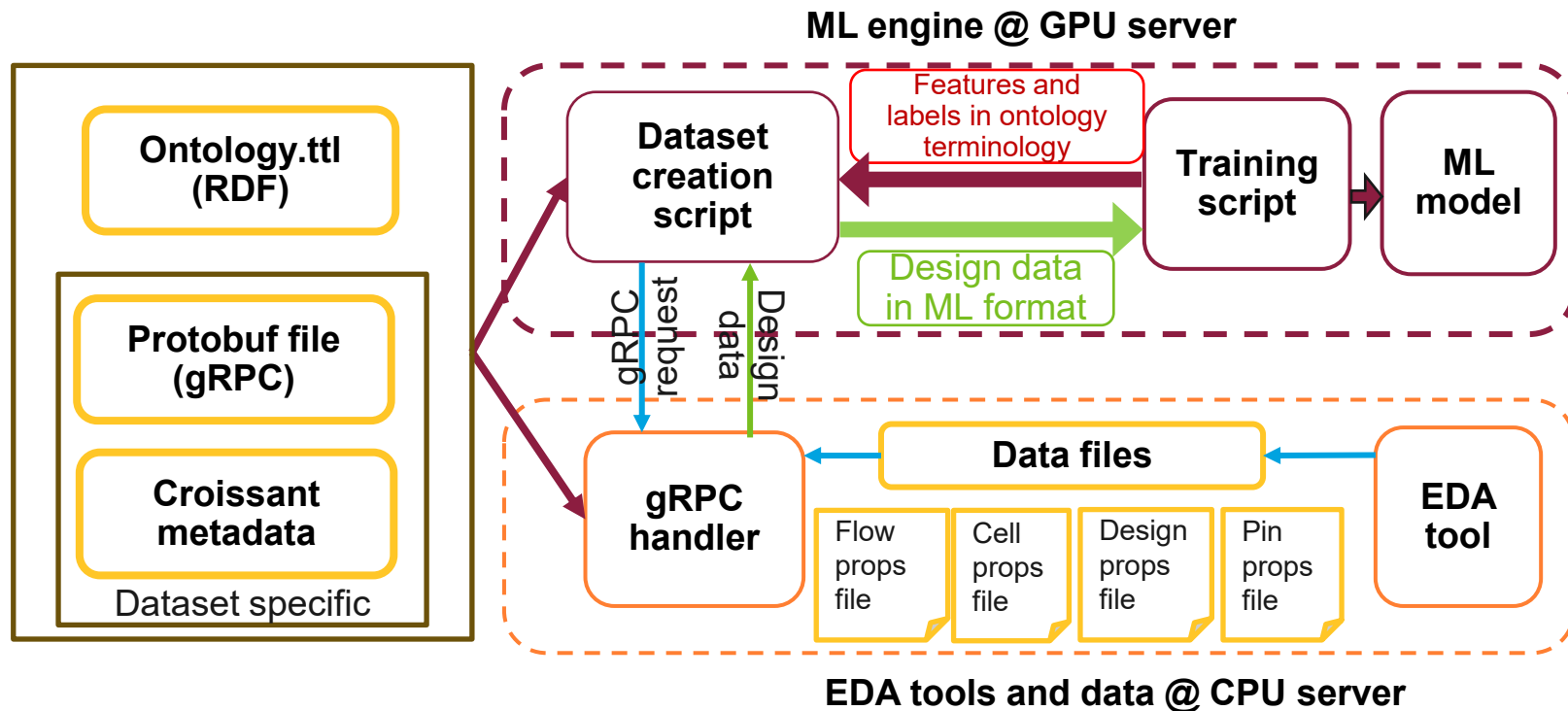
Enable AI-based EDA via:

- Standardized schema flows
- Ontology-based data relationships
- Tool-agnostic and dataset-agnostic data pipelines
- Seamless integration with EDA tools on CPU servers and ML on GPUs for training and inference



# Proposed Flow Overview

Slide credits: Vikram Gopalakrishnan, ASU



# Outline

- Challenges and motivation for ML EDA data pipelines
- Si2 open standards working group ML EDA data pipeline
  - Proposed flow overview
  - **Components of the flow with an example usecase**
- Conclusion



# Cell Arc Delay Example Usecase with CircuitOps as a Data Source

A model that predicts the cell delay with input slew and load capacitance.

- Features: Input pin slew, output load capacitance, cell type
- Labels: Arc delay

Required features and labels in ontology-compliant language

- Timing Graph with:
  - TimingEdge properties: edgeType, delayValue
  - sourceNode properties: worstSlew, libCellName
  - sinkNode properties: loadCapacitance

# Parts of the RDF Ontology for Usecase

Highlighting those classes and properties required for our usecase

## Owl Class Definitions

```
448 eda:Net a owl:Class ;
449   rdfs:label "Net" ;
450   eda:primaryId eda:netName ;
451   rdfs:comment "A collection of electrically connected pins and wires." .
452
453 eda:SignalNet a owl:Class ;
454   rdfs:subClassOf eda:Net ;
455   rdfs:label "Signal Net" .
456
457 eda:PowerNet a owl:Class ;
458   rdfs:subClassOf eda:Net ;
459   rdfs:label "Power Net" ;
460   rdfs:comment "Net that distributes power (e.g., VDD, VSS, or analog supplies)" .
461
462 eda:Wire a owl:Class ;
463   rdfs:label "Wire" ;
464   rdfs:comment "A physical segment of metal interconnect belonging to a net." .
465
466 eda:Via a owl:Class ;
467   rdfs:label "Via" ;
468   rdfs:comment "A vertical connection between different metal layers, belonging to a net." .
469
470 eda:Pin a owl:Class ;
471   rdfs:label "Pin" ;
472   eda:primaryId eda:pinName ;
473   rdfs:comment "A connection point on an instance or a port." .
474
475 eda:SignalPin a owl:Class ;
476   rdfs:subClassOf eda:Pin , [
477     a owl:Restriction ;
478     owl:onProperty eda:connectsTo ;
479     owl:allValuesFrom eda:SignalNet
480   ] ;
481   rdfs:comment "Pin that transmits logical or timing-related information." .
```

## Pin Properties

```
2002 ### Pin properties
2003 eda:pinName a owl:DatatypeProperty ;
2004   rdfs:domain eda:Pin ;
2005   rdfs:range xsd:string ;
2006   rdfs:label "pin name" .
2007
2008 eda:pinX a owl:DatatypeProperty ;
2009   rdfs:domain eda:Pin ;
2010   rdfs:range xsd:float ;
2011   rdfs:label "pin X" ;
2012   qudt:unit unit:MicroM .
2013
2014 eda:pinY a owl:DatatypeProperty ;
2015   rdfs:domain eda:Pin ;
2016   rdfs:range xsd:float ;
2017   rdfs:label "pin Y" ;
2018   qudt:unit unit:MicroM .
2019
2020 eda:worstSlew a owl:DatatypeProperty ;
2021   rdfs:domain eda:Pin ;
2022   rdfs:range xsd:float ;
2023   rdfs:label "worst slew" ;
2024   qudt:unit unit:NanoSecond .
2025
2026 eda:loadCapacitance a owl:DatatypeProperty ;
2027   rdfs:domain eda:Pin ;
2028   rdfs:range xsd:float ;
2029   rdfs:label "load capacitance" ;
2030   qudt:unit unit:FemtoFarad .
2031
2032 eda:isClockPin a owl:DatatypeProperty ;
2033   rdfs:domain eda:Pin ;
2034   rdfs:range xsd:boolean ;
2035   rdfs:label "is clock pin" .
```

# Parts of the RDF Ontology for Usecase

Highlighting those classes and properties required for the usecase

1268 **### Timing graph properties**  
1269 **eda:hasNode** a owl:ObjectProperty ;  
1270 **rdfs:domain** eda:TimingGraph ;  
1271 **rdfs:range** eda:TimingNode ;  
1272 **rdfs:label** "has node" ;  
1273 **rdfs:comment** "Relates a timing graph to a timing node contained within it." .  
1274  
1275 **eda:hasEdge** a owl:ObjectProperty ;  
1276 **rdfs:domain** eda:TimingGraph ;  
1277 **rdfs:range** eda:TimingEdge ;  
1278 **rdfs:label** "has edge" ;  
1279 **rdfs:comment** "Relates a timing graph to a timing edge contained within it." .  
1280  
1281 **eda:sourceNode** a owl:ObjectProperty ;  
1282 **rdfs:domain** eda:TimingEdge ;  
1283 **rdfs:range** eda:TimingNode ;  
1284 **eda:primaryId** eda:sourcePinName ;  
1285 **rdfs:label** "source node" ;  
1286 **rdfs:comment** "Relates a timing edge to its source node." .  
1287  
1288 **eda:sinkNode** a owl:ObjectProperty ;  
1289 **rdfs:domain** eda:TimingEdge ;  
1290 **rdfs:range** eda:TimingNode ;  
1291 **eda:primaryId** eda:sinkPinName ;  
1292 **rdfs:label** "sink node" ;  
1293 **rdfs:comment** "Relates a timing edge to its sink node." .  
1294  
1295 **eda:represents** a owl:ObjectProperty ;  
1296 **rdfs:domain** eda:TimingNode ;  
1297 **rdfs:range** [ a owl:Class ; owl:unionOf(eda:SignalPin eda:SignalNet eda:SignalPort) ] ;  
1298 **rdfs:label** "represents" ;  
1299 **rdfs:comment** "Links a timing node to the design element (pin, port, or net) it represents." .  
1300

## Timing graph

2534 **eda:delayValue** a owl:DatatypeProperty ;  
2535 **rdfs:domain** eda:TimingEdge ;  
2536 **rdfs:range** xsd:float ;  
2537 **rdfs:label** "delay value" ;  
2538 **qudt:unit** unit:NanoSecond .  
2539  
2540 **eda:edgeType** a owl:DatatypeProperty ;  
2541 **rdfs:domain** eda:TimingEdge ;  
2542 **rdfs:range** xsd:bool ;  
2543 **rdfs:label** "edge type" .  
2544  
2545 **eda:riseDelay** a owl:DatatypeProperty ;  
2546 **rdfs:domain** eda:TimingEdge ;  
2547 **rdfs:range** xsd:float ;  
2548 **rdfs:label** "rise delay" ;  
2549 **qudt:unit** unit:NanoSecond .  
2550  
2551 **eda:fallDelay** a owl:DatatypeProperty ;  
2552 **rdfs:domain** eda:TimingEdge ;  
2553 **rdfs:range** xsd:float ;  
2554 **rdfs:label** "fall delay" ;  
2555 **qudt:unit** unit:NanoSecond .  
2556

## Timing edge properties

Slide credits: Ivan Kissiov

# Croissant Metadata for CircuitOps

- Croissant is a metadata format for structured datasets, enabling interoperability and schema definition.
- Key elements:
  - *Dataset* – Represents the entire collection of files and metadata.
  - *FileObject* – Defines individual files in the dataset (e.g., CSVs, images, JSON).
  - *RecordSet* – Describes the structure of tabular data (e.g., design\_properties.csv).
  - *Field* – Specifies individual columns in a RecordSet, including data types and references.
- Reference: <https://docs.mlcommons.org/croissant/docs/croissant-spec.html#resources>

# Croissant Metadata for CircuitOps

## Dataset definition

```
1 {
2   "@type": "Dataset",
3   "@id": "source_dataset",
4   "name": "Source Dataset",
5   "description": "A dataset containing structured circuit design data.",
6   "distribution": [
7     {
8       "@type": "cr:FileObject",
9       "@id": "design_properties",
10      "name": "design_properties.csv",
11      "contentUrl": "data/design_properties.csv",
12      "encodingFormat": "text/csv"
13    },
14    {
15      "@type": "cr:FileObject",
16      "@id": "net_properties",
17      "name": "net_properties.csv",
18      "contentUrl": "data/net_properties.csv",
19      "encodingFormat": "text/csv"
20    },
21    {
22      "@type": "cr:FileObject",
23      "@id": "inst_properties",
24      "name": "inst_properties.csv",
25      "contentUrl": "data/inst_properties.csv",
26      "encodingFormat": "text/csv"
27    },
28    {
29      "@type": "cr:FileObject",
30      "@id": "pin_properties",
31      "name": "pin_properties.csv",
32      "contentUrl": "data/pin_properties.csv",
33      "encodingFormat": "text/csv"
34    },
35    {
36      "@type": "cr:FileObject",
37      "@id": "power_density",
38      "name": "power_density.png",
```

## RecordSet definition

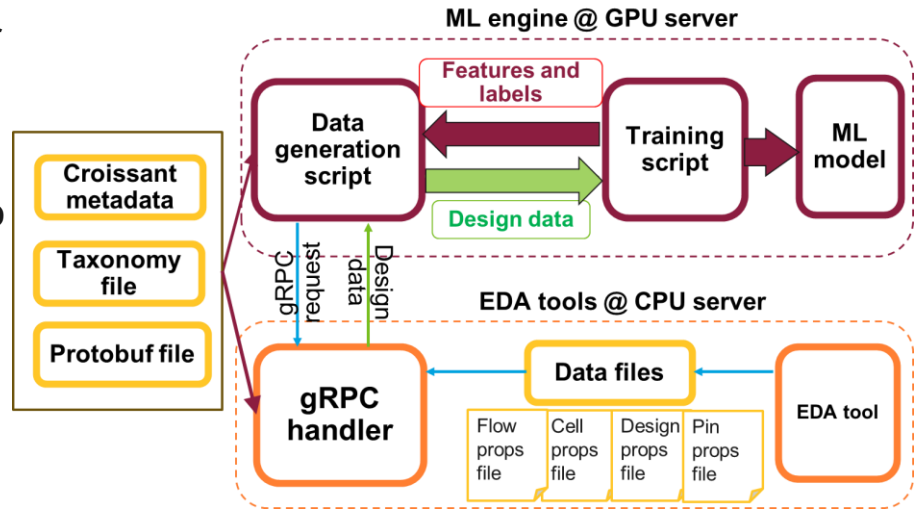
```
280 {
281   "@type": "cr:RecordSet",
282   "@id": "pin_properties",
283   "field": [
284     {
285       "@type": "cr:Field",
286       "@id": "pin_properties/pinName",
287       "name": "pinName",
288       "dataType": "sc:Text",
289       "source": {
290         "fileObject": { "@id": "pin_properties" },
291         "extract": { "column": "pin_name" }
292       }
293     },
294     {
295       "@type": "cr:Field",
296       "@id": "pin_properties/worstSlew",
297       "name": "worstSlew",
298       "dataType": "sc:Float",
299       "source": {
300         "fileObject": { "@id": "pin_properties" },
301         "extract": { "column": "pin_max_slew" }
302       }
303     },
304     {
305       "@type": "cr:Field",
306       "@id": "pin_properties/loadCapacitance",
307       "name": "loadCapacitance",
308       "dataType": "sc:Float",
309       "source": {
310         "fileObject": { "@id": "pin_properties" },
311         "extract": { "column": "load_capacitance" }
312       }
313     },
314     {
315       "@type": "cr:Field",
316       "@id": "pin_properties/inputCapacitance",
317       "name": "inputCapacitance",
318       "dataType": "sc:Float",
319       "source": {
320         "fileObject": { "@id": "pin_properties" },
321         "extract": { "column": "input_capacitance" }
```

# Conclusion and acknowledgment

- Using those components described above, we have developed automated scripts for the ML EDA data pipeline
- Utilizing existing ML and ontology infrastructure is the least resistant path to build data pipelines (Croissant, RDF, protobuf)
- **Standards require community input and commitment to adoption**

## Acknowledgments:

- Si2 open standard working group
- ASU VDA lab students



# Open-source ML EDA Research @ ASU

## Synthetic data generation techniques using AI

- Timing cones (MIMIC): <https://github.com/ASU-VDA-Lab/MIMIC>
- Layout heatmaps using diffusion models (DALI-PD): <https://github.com/ASU-VDA-Lab/DALI-PD>

## LLM chatbots and agents for OpenROAD

- EDA Corpus: Dataset and benchmarks to train chatbots for OpenROAD  
<https://github.com/OpenROAD-Assistant/EDA-Corpus>
- OpenROAD-Assistant: Chatbots QA and script generation  
<https://github.com/OpenROAD-Assistant/OpenROAD-Assistant>
- OpenROAD-Agent: Automated script generation and correction  
<https://github.com/OpenROAD-Assistant/OpenROAD-Agent>