# SwSS (entry-point), deep dive

Understanding, how the SwSS container works

# Objectives

- Learn the SWSS container input options
- Why use ZMQ and Redis at the same time?
- Learn to configure input arguments
- Learn the initializations (SAI and SAI-Redis)
- Learn the DB connections
- Learn about SAI Attributes
- Learn the program flow

XFLOW
RESEARCH

# main.cpp

# Argument parsing and config variables

- command line options
- for container initialization

| | |
|---|---|
| b | batch size |
| i | asic instance string input |
| m | MAC-Address |
| r | enable/disable recording inputs |
| d | writing directory |
| h | help |
| s | enable synchronization |
| z | sai_deserialize_redis_communication_mode |
| f | file name |
| k | max bulk size limit in bulk mode |
| q | zmq server address |

XFLOW
RESEARCH

# Why use ZMQ and Redis at the same time?

| Redis | ZMQ (Zero M Queue) |
|---|---|
| **Purpose (use case):** | |
| - Data Store<br>    - Key Value pairs<br>- Benefits<br>    - **Language independent** Interface<br>    - **persistent** storage<br>    - **replication** and i**nter-process** communication | - **RTPS**: Real TIme Publish Subscribe, library<br>    - No specific message format<br>    - **low latency**<br>- complex routing and filtering<br>- **Distributed systems** / environments<br>- **custom protocols**<br>- Different operation modes (**IPC**, **TCP** and **Memory**) |
| **Message Types:** | |
| - Pub / Sub | - Pub / Subs<br>- Request / Reply<br>- Push / Pull<br>- Pipeline (Buffer/Queue) |
| **Storage Type:** | |
| - In-memory | |

**XFLOW** RESEARCH

# Why ZMQ in SONiC and not any other RTPS or DDS

1) Telemetry and Monitoring
   - sonic **generates** a lot of **telemetry data**
     - network performance
     - traffic statistics
     - device health metrics (etc)
   - stream and deliver to visualization applications
1) Orchestration and Automation
   - Enabling network managers to **communicate** with **external** and **3rd party components** for **automation purposes**.
1) Example use case is,
   - Defining the **API**s for **remote configuration** of the containers or the whole sonic environment. e.g in g**NMI** (g**RPC**, Network Management Interface) the commands from client get to the server and from there they are **communicated to SONiC** architecture through **ZMQ**

XFLOW
RESEARCH

# What is a Recorder

- Interface class to use recorder instances
- Affiliated to logging

```cpp
namespace swss {
};

/* Interface to access recorder classes */
class Recorder {
public:
    static Recorder& Instance();
    static const std::string DEFAULT_DIR;
    static const std::string REC_START;
    static const std::string SWSS_FNAME;
    static const std::string SAIREDIS_FNAME;
    static const std::string RESPPUB_FNAME;

    Recorder() = default;
    /* Individual Handlers */
    SwSSRec swss;
    SaiRedisRec sairedis;
    ResPubRec respub;
};
```

XFLOW
RESEARCH

# Recorder Instance

- on input options, decision made if recorder has to be enabled or left disabled
  - swss,
  - respub,
  - sairedis
- Parameters set on initialization
  - setLocation
  - setFileName
  - startRec

```
/* Initialize sairedis recording parameters */
Recorder::Instance().sairedis.setRecord(
    (record_type & SAIREDIS_RECORD_ENABLE) == SAIREDIS_RECORD_ENABLE
);
Recorder::Instance().sairedis.setLocation(record_location);
Recorder::Instance().sairedis.setFileName(sairedis_rec_filename);

/* Initialize sairedis */
initSaiApi();
initSaiRedis();

/* Initialize remaining recorder parameters  */
Recorder::Instance().swss.setRecord(
    (record_type & SWSS_RECORD_ENABLE) == SWSS_RECORD_ENABLE
);
Recorder::Instance().swss.setLocation(record_location);
Recorder::Instance().swss.setFileName(swss_rec_filename);
Recorder::Instance().swss.startRec(true);

Recorder::Instance().respub.setRecord(
    (record_type & RESPONSE_PUBLISHER_RECORD_ENABLE) ==
    RESPONSE_PUBLISHER_RECORD_ENABLE
);
Recorder::Instance().respub.setLocation(record_location);
Recorder::Instance().respub.setFileName(responsepublisher_rec_filename);
Recorder::Instance().respub.startRec(false);
```
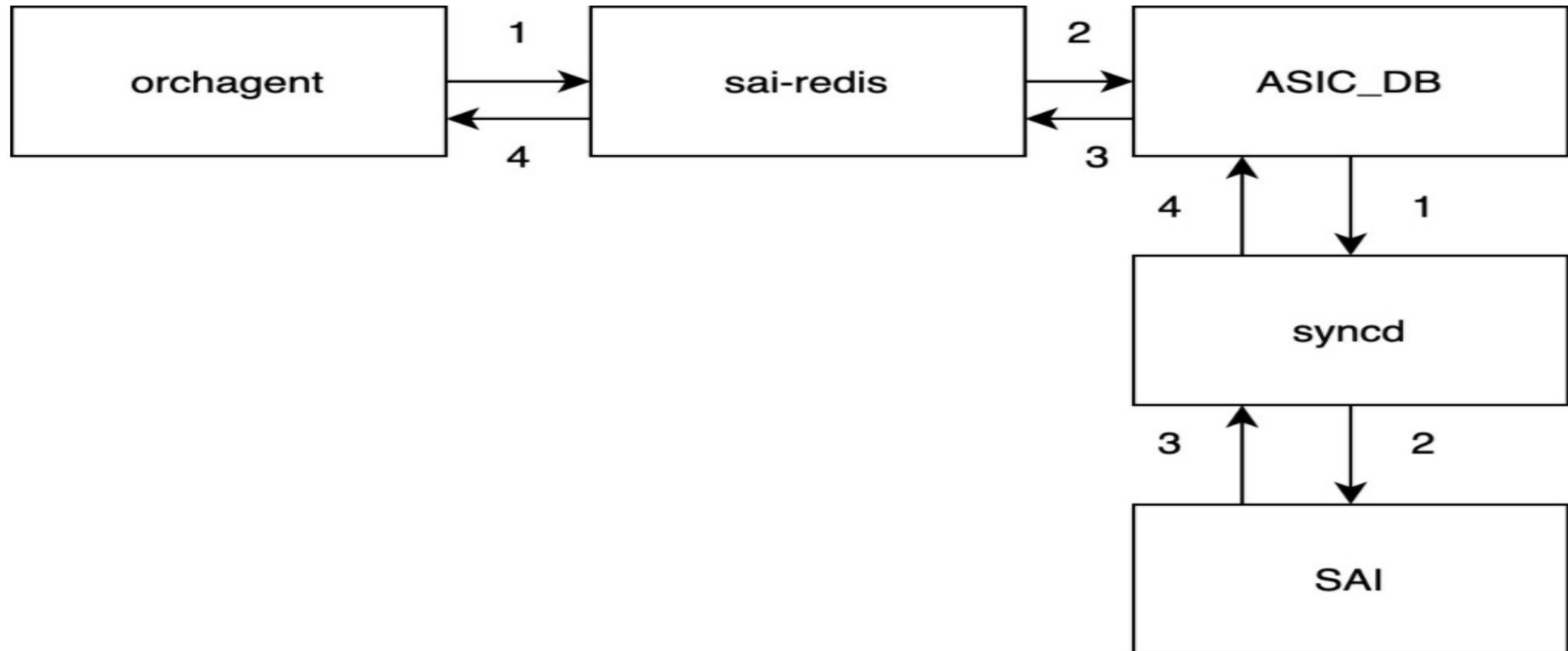
XFLOW
RESEARCH

# Initialize SAI

Initialize SAI API
Initialize SAI Redis

```
/* Initialize sairedis */
initSaiApi();
initSaiRedis();
```

# Orch-agent to ASICs: Information flow

# Instantiate DB

- DB initializations
  - APPL_DB
  - CONFIG_DB
  - State_DB

```
320    int main(int argc, char **argv)
465
466        // Instantiate database connectors
467        DBConnector appl_db("APPL_DB", 0);
468        DBConnector config_db("CONFIG_DB", 0);
469        DBConnector state_db("STATE_DB", 0);
470
```

# Initialize ZMQ communication (Conditional)

```
470
471        // Instantiate ZMQ server
472        shared_ptr<ZmqServer> zmq_server = nullptr;
473        if (enable_zmq)
474        {
475            SWSS_LOG_NOTICE("Instantiate ZMQ server : %s", zmq_server_address.c_str());
476            zmq_server = make_shared<ZmqServer>(zmq_server_address.c_str());
477        }
478        else
479        {
480            SWSS_LOG_NOTICE("ZMQ disabled");
481        }
482
```

- if **zmq_server** is specified in **option** then this **if command executes** and initializes zmq_communication.

XFLOW
RESEARCH

# switch_type and sai_attributes

- global variable

```
65      string gMySwitchType = "";
```

| - | dpu | data processing unit |
|---|-----|---------------------|
| - | voq | virtual output queue |
| - | fabric | communication medium within multiple components inside a device |
| - | chassis-packet | end-to-end information flow |

- check switch type and then configuring the sai_attribute_t object i.e attr
- vector<sai_attribute_t> attrs

```
483     // Get switch_type
484     getCfgSwitchType(&config_db, gMySwitchType);
485
486     sai_attribute_t attr;
487     vector<sai_attribute_t> attrs;
488
489     attr.id = SAI_SWITCH_ATTR_INIT_SWITCH;
490     attr.value.booldata = true;
491     attrs.push_back(attr);
492
493     if (gMySwitchType != "dpu")
494     {
495         attr.id = SAI_SWITCH_ATTR_FDB_EVENT_NOTIFY;
496         attr.value.ptr = (void *)on_fdb_event;
497         attrs.push_back(attr);
498     }
```

XFLOW
RESEARCH

# sai_attributes :    atrr -> (id, value(booldata, ptr) )

| Sr | SAI Attribute IDs | IF conditional on global variables | Reference Image |
|---|---|---|---|
| 1. | SAI_SWITCH_ATTR_INIT_SWITCH | | `489    attr.id = SAI_SWITCH_ATTR_INIT_SWITCH;`<br>`490    attr.value.booldata = true;`<br>`491    attrs.push_back(attr);` |
| 2. | SAI_SWITCH_ATTR_FDB_EVENT_NOTIFY | gSwitchState != "**dpu**" | `495    attr.id = SAI_SWITCH_ATTR_FDB_EVENT_NOTIFY;`<br>`496    attr.value.ptr = (void *)on_fdb_event;`<br>`497    attrs.push_back(attr);` |
| 3. | SAI_SWITCH_ATTR_PORT_STATE_CHANGE_NOTIFY | | `500    attr.id = SAI_SWITCH_ATTR_PORT_STATE_CHANGE_NOTIFY;`<br>`501    attr.value.ptr = (void *)on_port_state_change;`<br>`502    attrs.push_back(attr);` |
| 4. | SAI_SWITCH_ATTR_SHUTDOWN_REQUEST_NOTIFY | | `504    attr.id = SAI_SWITCH_ATTR_SHUTDOWN_REQUEST_NOTIFY;`<br>`505    attr.value.ptr = (void *)on_switch_shutdown_request;`<br>`506    attrs.push_back(attr);` |
| 5. | SAI_SWITCH_ATTR_PORT_HOST_TX_READY_NOTIFY | | `508    attr.id = SAI_SWITCH_ATTR_PORT_HOST_TX_READY_NOTIFY;`<br>`509    attr.value.ptr = (void *)on_port_host_tx_ready;`<br>`510    attrs.push_back(attr);` |
| 6. | SAI_SWITCH_ATTR_SRC_MAC_ADDRESS | gSwitchState != "**fabric**" **&&** getMACAddress | `514    attr.id = SAI_SWITCH_ATTR_SRC_MAC_ADDRESS;`<br>`515    memcpy(attr.value.mac, gMacAddress.getMac(), 6);`<br>`516    attrs.push_back(attr);` |

# Redis Communication Mode at Initialization

- global variable

```
55   bool gSyncMode = false;
```

- default **false**
- set to **true** on **-s** option in cli. i.e **synchronization**
- **deprecated,** use **-z** instead **(deserialize)**

```
391        case 's':
392            gSyncMode = true;
393            SWSS_LOG_NOTICE("Enabling synchronous mode");
394        break;
```

```
522    if (gSyncMode)
523    {
524        SWSS_LOG_WARN("sync mode is depreacated, use -z param");
525
526        gRedisCommunicationMode = SAI_REDIS_COMMUNICATION_MODE_REDIS_SYNC;
527    }
```

XFLOW
RESEARCH

# OrchDaemon (init() and start() called in main.cpp)

```cpp
320    int main(int argc, char **argv)
742        shared_ptr<OrchDaemon> orchDaemon;
743        if (gMySwitchType != "fabric")
744        {
745            orchDaemon = make_shared<OrchDaemon>(&appl_db, &config_db, &state_db, chassis_app_db.get(), zmq_server.get());
746            if (gMySwitchType == "voq")
747            {
748                orchDaemon->setFabricEnabled(true);
749                orchDaemon->setFabricPortStatEnabled(true);
750                orchDaemon->setFabricQueueStatEnabled(false);
751            }
752        }
753        else
754        {
755            orchDaemon = make_shared<FabricOrchDaemon>(&appl_db, &config_db, &state_db, chassis_app_db.get(), zmq_server.get());
756        }
757
758        if (!orchDaemon->init())
759        {
760            SWSS_LOG_ERROR("Failed to initialize orchestration daemon");
761            exit(EXIT_FAILURE);
762        }
```

```cpp
772
773            orchDaemon->start();
774
```

XFLOW RESEARCH

# orchdaemon attributes

4 databases

- appl
- config
- state
- chassisApp

1 zmq server

```
orchagent > C orchdaemon.h > ...
 57    class OrchDaemon
 85    private:
 86        DBConnector *m_applDb;
 87        DBConnector *m_configDb;
 88        DBConnector *m_stateDb;
 89        DBConnector *m_chassisAppDb;
 90        ZmqServer *m_zmqServer;
 91
 92        bool m_fabricEnabled = false;
 93        bool m_fabricPortStatEnabled = true;
 94        bool m_fabricQueueStatEnabled = true;
 95
 96        std::vector<Orch *> m_orchList;
 97        Select *m_select;
 98
 99        std::chrono::time_point<std::chrono::hi
```

```
 57    class OrchDaemon
 58    {
 59    public:
 60        OrchDaemon(DBConnector *, DBConnector *, DBConnector *, DBConnector *, ZmqServer *);
 61        ~OrchDaemon();
```

XFLOW
RESEARCH

# Child / Sub-Orch

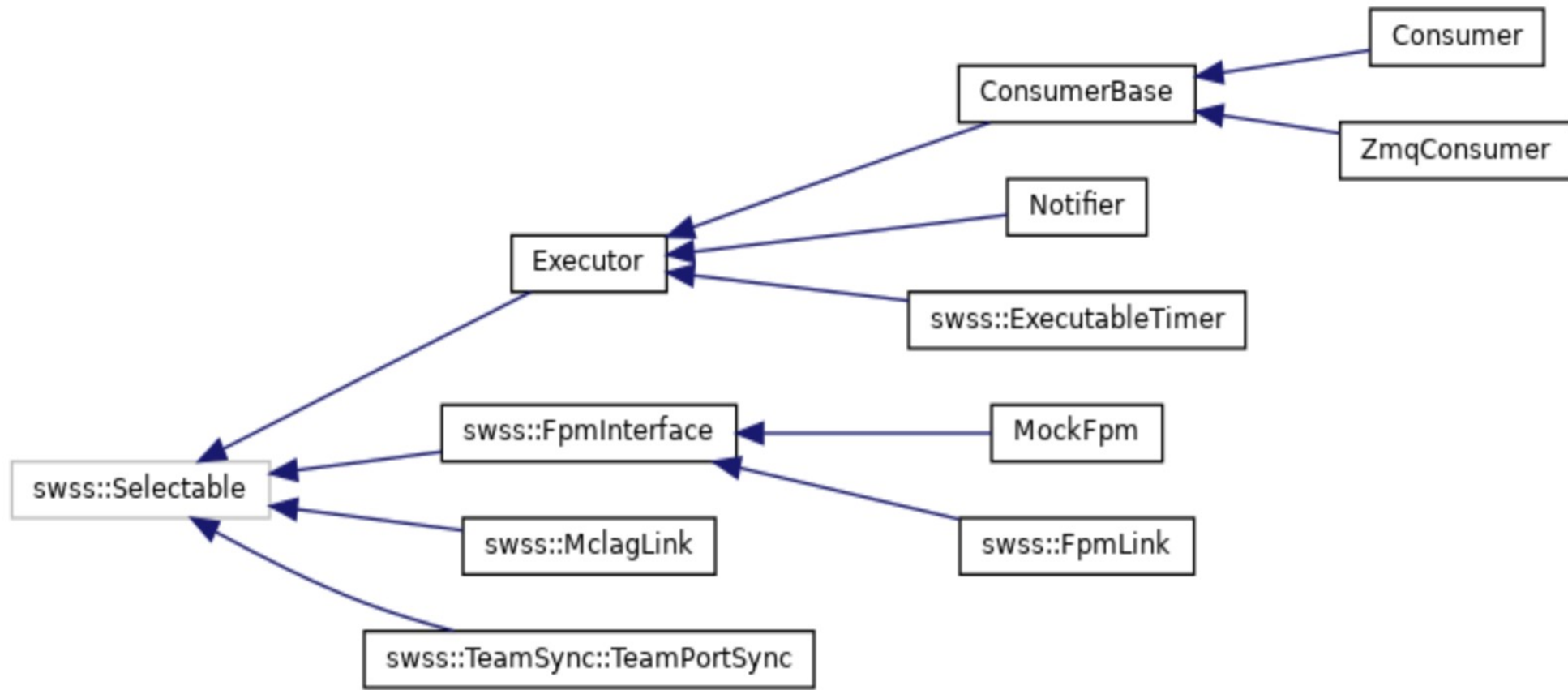Each of which inherits the **orchdaemon** class.

```
30  /*
31   * Global orch daemon variables
32   */
33  PortsOrch *gPortsOrch;
34  FabricPortsOrch *gFabricPortsOrch;
35  FdbOrch *gFdbOrch;
36  IntfsOrch *gIntfsOrch;
37  NeighOrch *gNeighOrch;
38  RouteOrch *gRouteOrch;
39  NhgOrch *gNhgOrch;
40  NhgMapOrch *gNhgMapOrch;
41  CbfNhgOrch *gCbfNhgOrch;
42  FgNhgOrch *gFgNhgOrch;
43  AclOrch *gAclOrch;
44  PbhOrch *gPbhOrch;
45  MirrorOrch *gMirrorOrch;
46  CrmOrch *gCrmOrch;
47  BufferOrch *gBufferOrch;
48  QosOrch *gQosOrch;
49  SwitchOrch *gSwitchOrch;
50  Directory<Orch*> gDirectory;
51  NatOrch *gNatOrch;
52  PolicerOrch *gPolicerOrch;
53  MlagOrch *gMlagOrch;
54  IsoGrpOrch *gIsoGrpOrch;
55  MACsecOrch *gMacsecOrch;
56  CoppOrch *gCoppOrch;
57  P4Orch *gP4Orch;
58  BfdOrch *gBfdOrch;
59  Srv6Orch *gSrv6Orch;
60  FlowCounterRouteOrch *gFlowCounterRouteOrch;
61  DebugCounterOrch *gDebugCounterOrch;
62  MonitorOrch *gMonitorOrch;
63
```
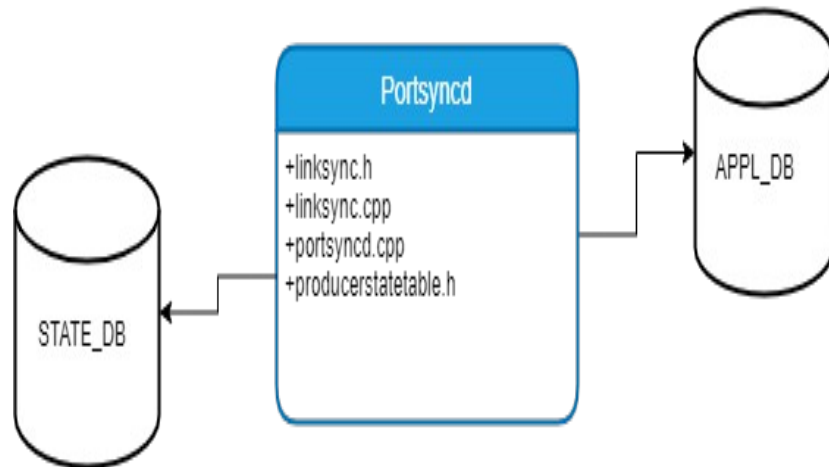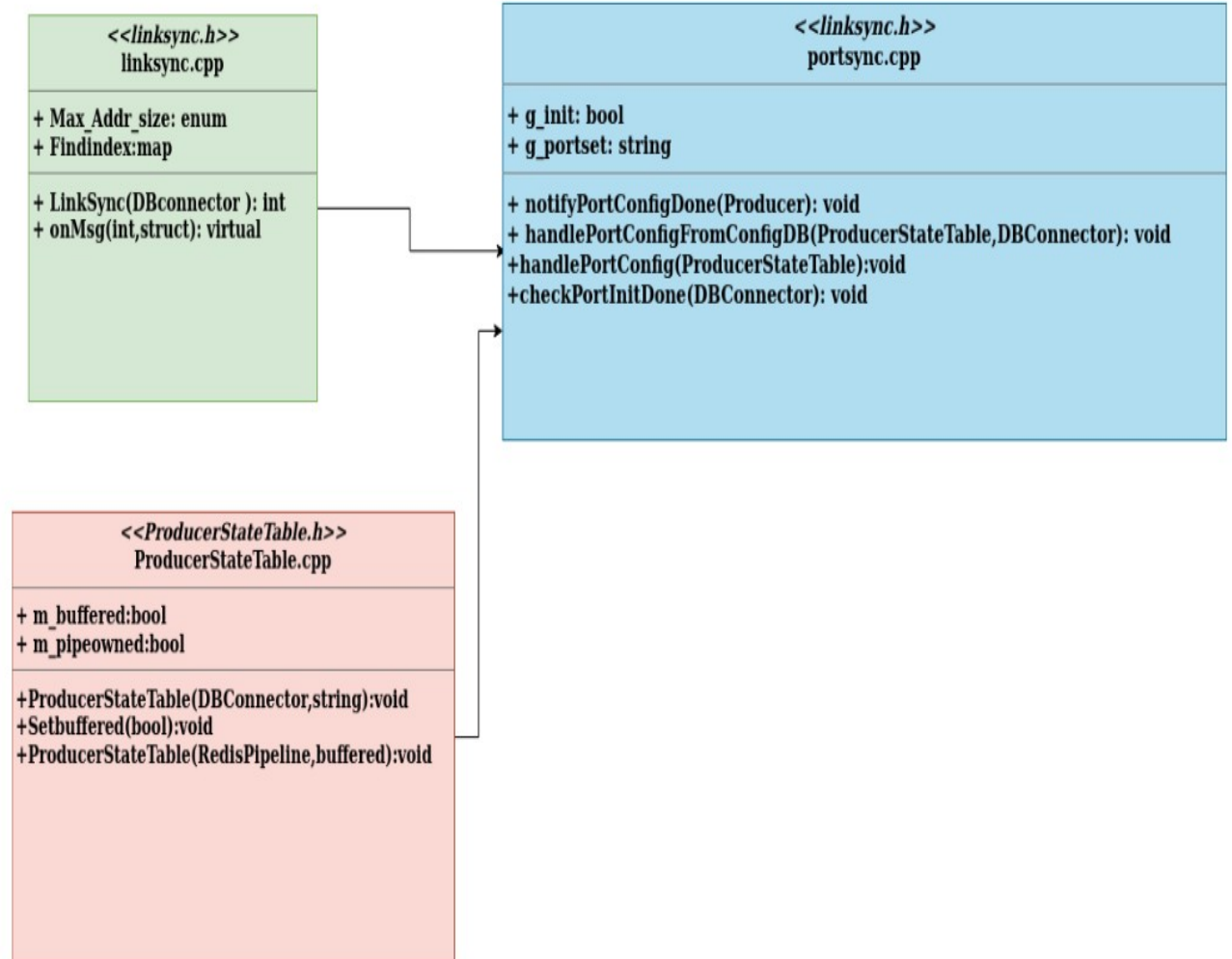
# SwSS repo

SwSS Repo
- orchagent
- portsyncd
- teamsyncd
- neighsyncd (lldp_syncd)
- swssconfig
- fpmsyncd
- natsyncd
- mclagsyncd
- gearsyncd

# orchagent.cpp

# SwSS Class Diagram

# **portsyncd** class diagram

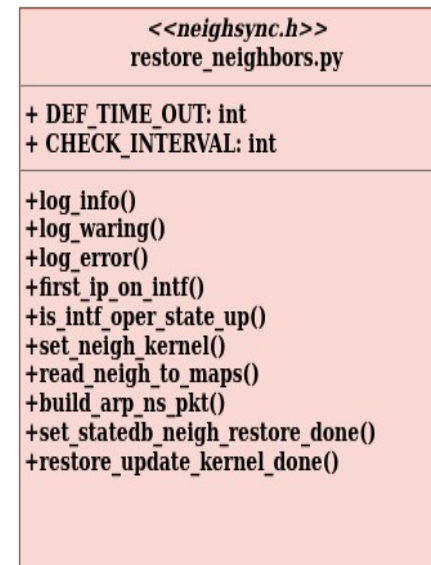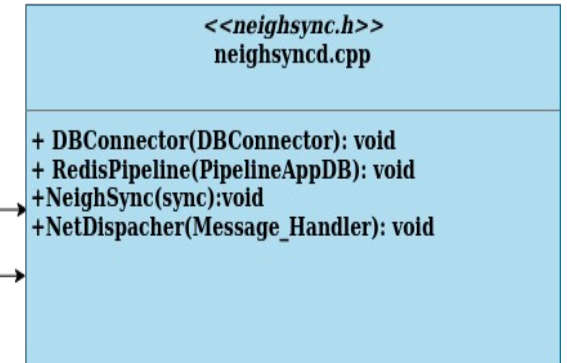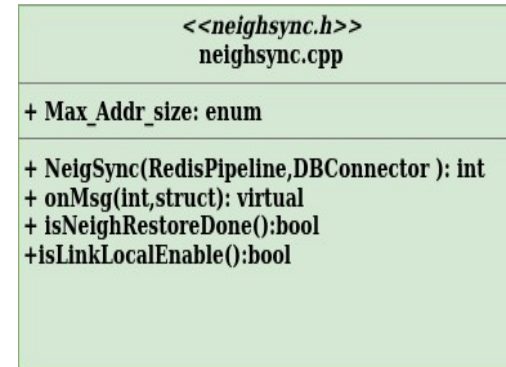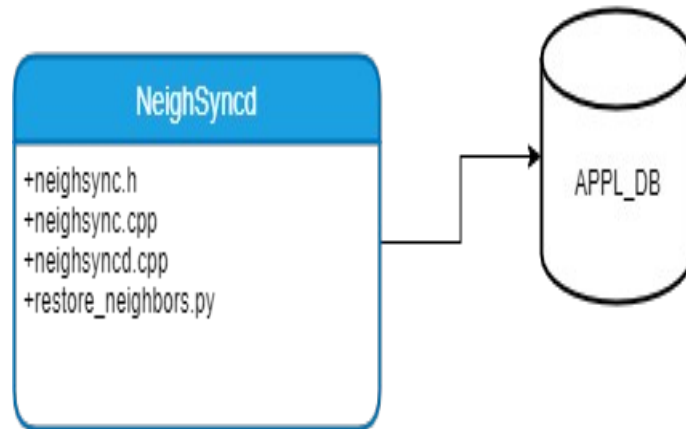# portsyncd **responsibilities**

- physical port information and synchronization with database
- Contains, **ProducerStateTable** which manages records of all **available ports** and their **states**
- **producerstatetable.h** is used by **linksync** file to manage the synchronization of ports with Database
- both these files are used by **portsyncd.cpp** to manage and notify port speed, lanes and mtu information to **APPL_DB** and **STATE_DB**
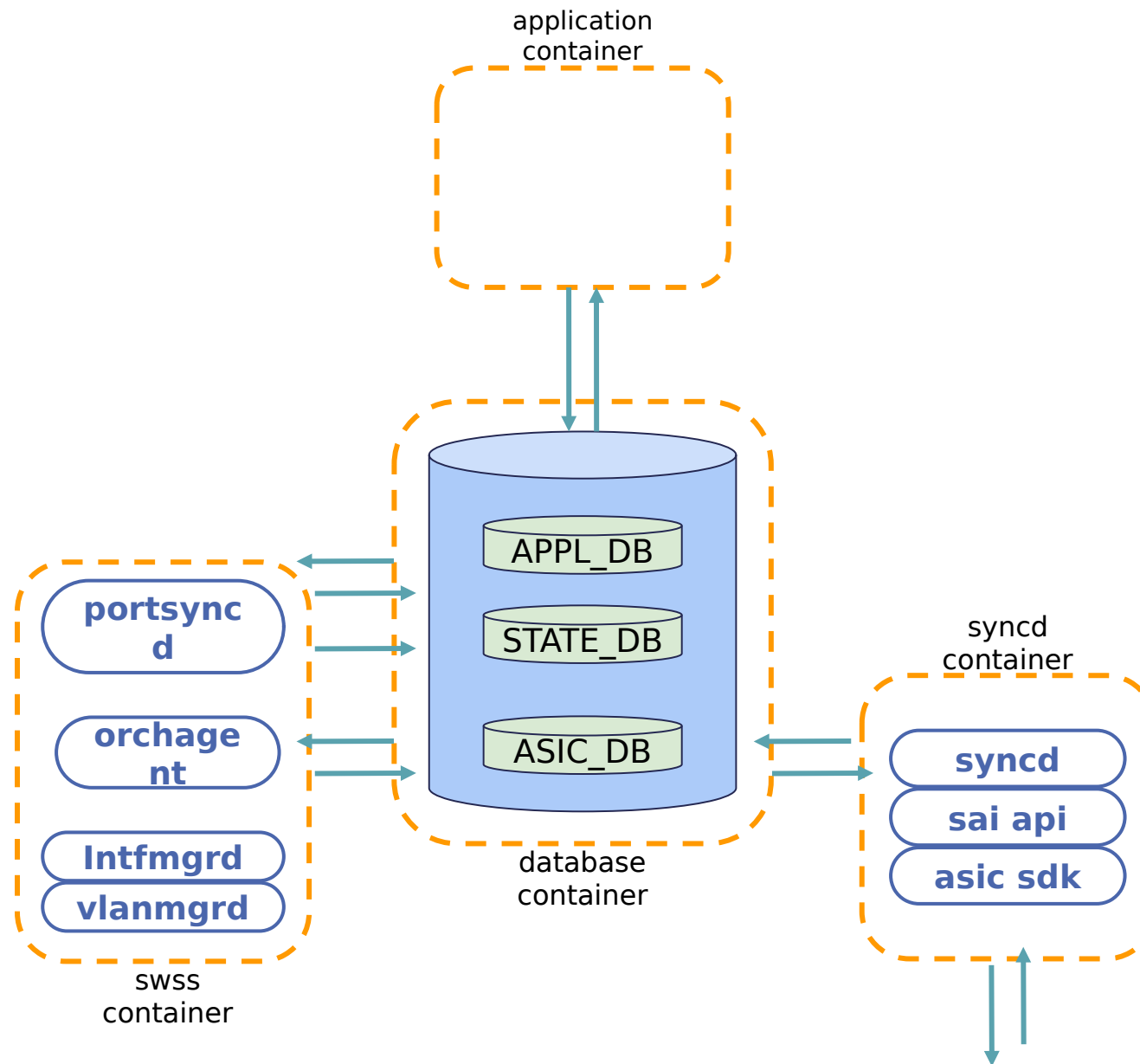
# **neighsyncd** class diagram

**<<neighsync.h>>**
neighsync.cpp

+ Max_Addr_size: enum

+ NeigSync(RedisPipeline,DBConnector ): int
+ onMsg(int,struct): virtual
+ isNeighRestoreDone():bool
+isLinkLocalEnable():bool

**<<neighsync.h>>**
neighsyncd.cpp

+ DBConnector(DBConnector): void
+ RedisPipeline(PipelineAppDB): void
+NeighSync(sync):void
+NetDispacher(Message_Handler): void

Overview of NeighSyncd

NeighSyncd

+neighsync.h
+neighsync.cpp
+neighsyncd.cpp
+restore_neighbors.py

APPL_DB

**<<neighsync.h>>**
restore_neighbors.py

+ DEF_TIME_OUT: int
+ CHECK_INTERVAL: int

+log_info()
+log_waring()
+log_error()
+first_ip_on_intf()
+is_intf_oper_state_up()
+set_neigh_kernel()
+read_neigh_to_maps()
+build_arp_ns_pkt()
+set_statedb_neigh_restore_done()
+restore_update_kernel_done()

XFLOW
RESEARCH

# neighsyncd **responsibilities**

- responsible for neighbour related activities like record management in data plane for L2 purposes.
- **neighsync.cpp** is there to discover and synchronize all the neighbors and maintain a **neighbors_table**
- Once all neighbors are synchronized it uses the **neighsyncd.cpp** file. In case of warm start the neighbor_table is read and cached onto the hashmap and keeps it updated from there on
- **restore_neighbors.py** is used for restoring the neighbor table from the **kernel** during system **warm reboot**. Agent **supervisord** in SwSS gets started on docker container startup. on warm reboot enabled it sets **stateDB flag** so neighsyncd can continue the reconciliation process.

application container

portsync d

orchage nt

Intfmgrd

vlanmgrd

swss container

APPL_DB

STATE_DB

ASIC_DB

database container

syncd container

syncd

sai api

asic sdk

28

# Message Broker

| Container **Process** | **Database** Container |
|:---:|:---:|
| **\*syncd** ||

| **Container** | **Container** |
|:---:|:---:|
| **Redis_DB (**Container**)** ||

| **Application** Container | **ASIC** hardware |
|:---:|:---:|
| **SwSS** ||

XFLOW RESEARCH

# Kernel Communication

| Application Container | | SwSS Container | |
|---|---|---|---|
| **Receive** from kernel | **Trans** to kernel | **Receive** from kernel /DB | **Trans** to kernel |
| **Daemon** process from **DB** (pushed by *syncd - SwSS) | **_syncd** process | __**synd** process (writes to DB and *mgrd - SwSS proceeds) | __**mngrd** process |

XFLOW RESEARCH

# Key Terminologies

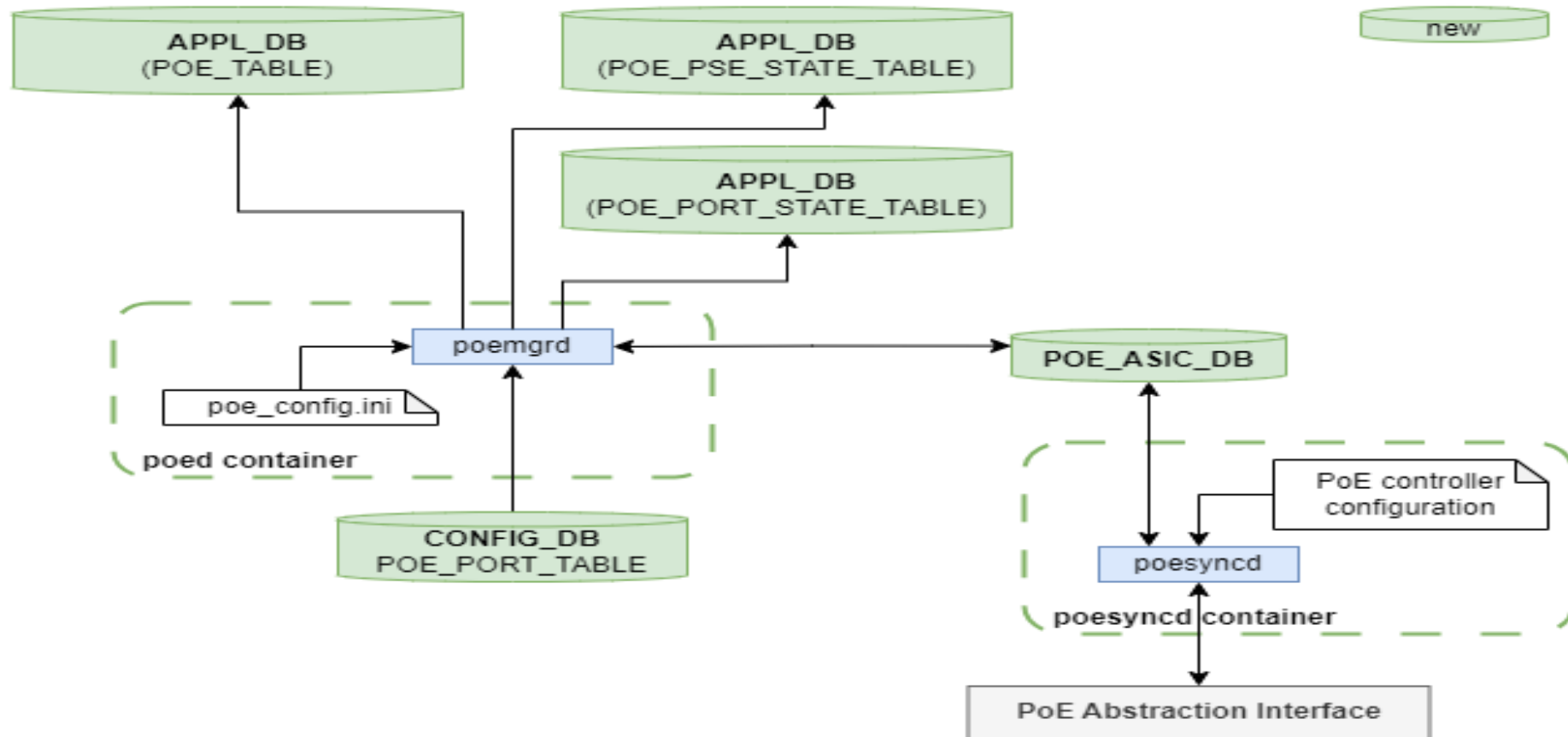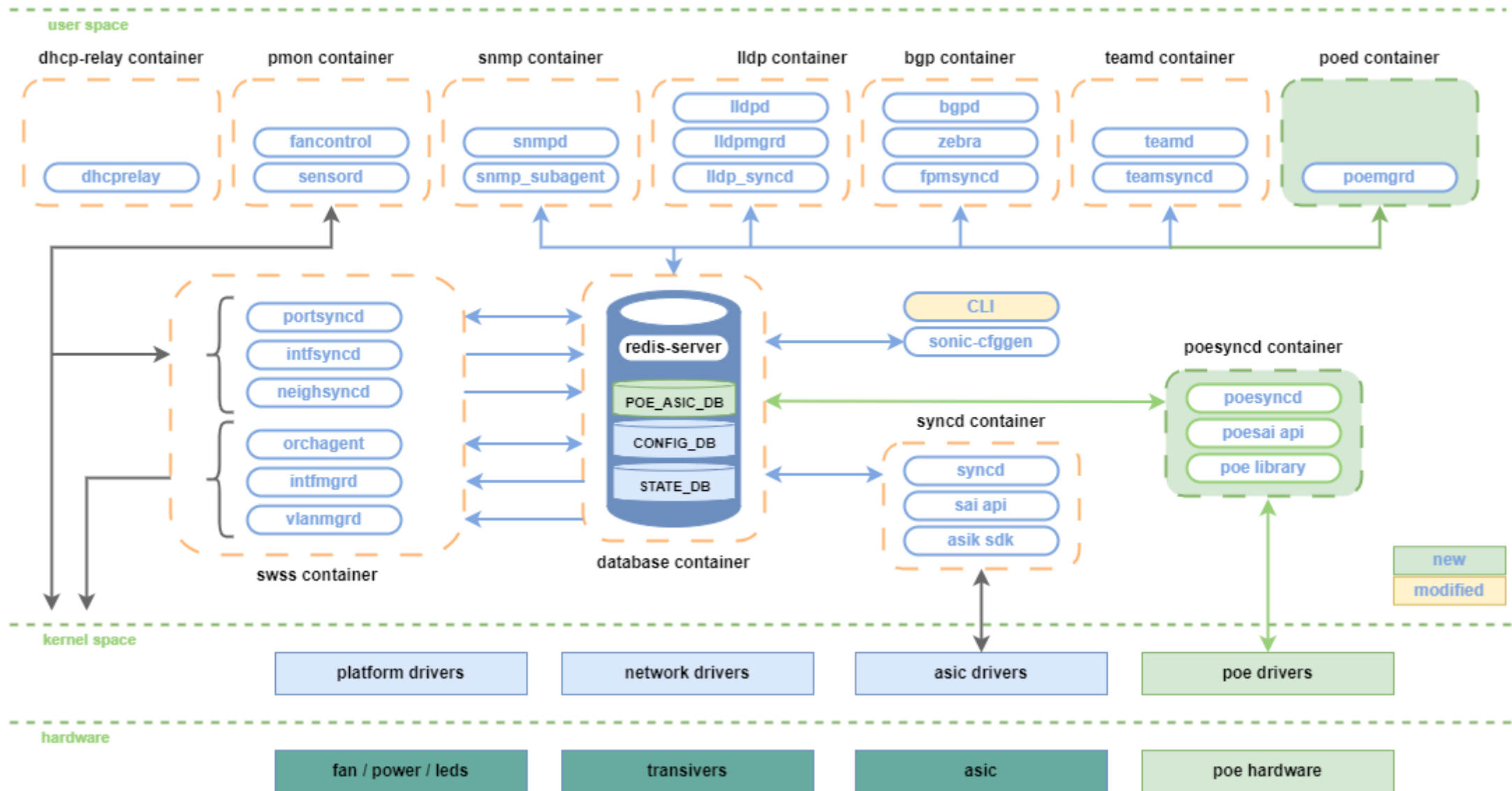| | |
|---|---|
| **Application Container** | Each application container,<br>1) Reads messages from kernel (netlink or socket or TCP)<br>2) Writes onto APPL_DB for other containers, through **syncd subprocesses** |
| **Daemon** | - Receives **netlink events** from the **kernel space**<br>- **Shares back** the responses in **few cases**. |
| **Syncd** | - **Publish** to **DB**s, the current **state** of the **container** with the environment |
| **Mgrd** | - **Reads APPL_DB** and **sends back** the **response** from **SwSS** to the OS **kernel**.<br>- **Example**:  config interface in linux kernel etc |
| **Orch** | - **Reads** from **APPL_DB, CONFIG_DB and STATE_DB**<br>- **Configure** the **AISCS** on start-up through **switch attributes**<br>- **Dictates** to **ASICS**, south-bound forwarding devices |
| **Cont -> Cont Communication** | **Principle**:<br>- The SONiC architecture follows **publish-subscribe** architecture.<br>- **containers** only **interact** through the **DB** and **not directly** in any way. |
| **South-bound communication**<br>Control-plane -> ASIC | 1) orch-subprocess at SwSS<br>2) writes to ASIC_DB, through sairedisapi<br>3) Eventhandler at sync_d container receives and pasess to the ASIC |
| **North-bound communication**<br>Application Cont -> SwSS | **1) Event** from **ASICS** through **driver** (kernel space), written to **ASIC_DB** by **syncd**<br>**2) Orchagent**, reads the state and **notifies APPL_DB** (notify apps) and **ASIC_DB** (notify kernel)<br>3) reply comes back from the kernel to **\*syncd.** e.g **portsyncd** |

XFLOW
RESEARCH

# PoE proposed architecture
## By NVIDIA

XFLOW
RESEARCH

# PoE: Power over Ethernet (Generally)

# Switch Type

| Sr # | Terminology | Details | Additional Reference |
|---|---|---|---|
| 1 - | **dpu**<br>**(**data processing unit**)** | - Involves tasks like switching and forwarding | - Another important aspect is that **DASH** (Disaggregated APIs for Sonic Hosts) **enabled DPU**s can be **used** in **SONiC infrastructure** out of the box. |
| 2 - | **voq**<br>**(**virtual output queue**)** | - Used to **control** / limit the flow of information on egress port as per network device capability | |
| 3 - | **fabric**<br>**(**communication within multiple components inside a device**)** | - Defines the **medium** and **protocol** on which **two** different **hardware** boards communicate (within one larger **system** or abstraction) e.g **cross GPU fabric** (common in **discussions** related to **hardware** in **AI-workgroup**) | Three types of switching fabrics<br> |
| 4 - | **chassis / modular**<br>**(**dedicated rack type switches with extended functionalities**)** | - Rather than to **group** multiple switches through cables (**stackable** switch), we can opt for a chassis switch<br>- switch **cards** can be swapped as per requirements<br>- Most **flexibility** |  |

XFLOW RESEARCH

# DASH (Disaggregated APIs for SONiC Hosts)

- **Another project** by **Microsoft**
- Implements the **logic** for **hardware reconfigurability** through **software**
- Recently incorporated in SONiC
- **Extends** functionality of **SAI** for **DPU**s or **Smart NIC**s
- DASH enabled DPUs are compatible with SONiC deployed datacenter / infrastructure

XFLOW
RESEARCH

# SAI player (SwSS)

- It records every sai operation/command at every line
- And **saves** as a **recording file**
- Can be **played back** for debugging purposes
- Only **specific** to a certain **ASIC**
- Playback **not** applicable across **different ASIC**
- **Move** recording **file** to **syncd** and **replay**

change. Also it must also be used on the same ASIC with the same SAI version. Cross ASIC replaying is not supported.

Here is possible scenario:

1. bug is spotted, ASIC is not configured as expected
2. Remove comment lines from the recording file such as the ones which begin with `|#|`
3. copy recording file `sairedis.2017-04-27.02:47:15.674566.rec` from `swss` docker to `syncd` docker

```
docker cp swss:/sairedis.2017-04-27.02:47:15.674566.rec .
docker cp ./sairedis.2017-04-27.02:47:15.674566.rec syncd:/
```

4. stop all sonic processes and clean the redis DB for fresh start

```
docker exec -it swss killall orchagent
redis-cli flushall
docker exec -it syncd killall syncd
docker exec -id syncd service syncd start
```

5. replay recording in `syncd` docker

```
docker exec -it syncd saiplayer sairedis.2017-04-27.02:47:15
```

6. confirm that ASIC is in bad state as found in the first place
7. describe what the problem is and send recoding file to vendor for investigation

XFLOW
RESEARCH

# Any Questions

XFLOW
RESEARCH

1) what is gear in swss? HLD was mentioned in a new proposal of PoE HLD community meeting. (It is a microservices open source framework incorporated in SONiC)
2) **fpmsyncd**, **teamsyncd, lldp_syncd** are exceptions. They are a **part of Application containers** but **code exists in SwSS** container as well. **Why**?
3) In the document shared about SwSS, what is the use of **restoreneighbors.py.** And how does this routine of **restore neighbor table** from **kernel** work in event of **warm reboot**?
4) What is the **broker** component in
   a) **across**-**containers** communication?
   b) **demon** and **database** communication?
   c) **database** and **ASIC** (Application plane and Control plane) communication?
      i) Is it bi-directional?
5) Do core devices also use only 4 wires(2 pairs) over and RJ-45 connector? There would be OFN, right? Is there a limit of how much bandwidth gets transmitted on copper wires?

# Khappa Scene kab karna hai ...........