

# In-Network Computing for App-Centric Micro-Services

draft-sarathchandra-coin-appcentres-03

D. Trossen, C. Sarathchandra, M. Boniface

COIN RG IETF 109 @ 18.11.2020

# Premise of the Draft

“The application-centric deployment of 'Internet' services has increased over the past ten years with many million applications providing user-centric services, executed on increasingly more powerful smartphones that are supported by Internet-based cloud services in distributed data centres, the latter mainly provided by large scale players such as Google, Amazon and alike. *This draft outlines a vision of evolving those data centres towards executing app-centric micro-services; we dub this evolved data centre as an **AppCentre**.*”

- Draft outlines use cases and research challenges for this vision

# General Structure

1 Introduction .....	4
2 Terminology .....	5
3 Use Cases .....	5
3.1 Mobile Application Function Offloading .....	5
3.2 Interactive Real-time Applications .....	7
3.3 Distributed AI .....	7
3.4 Content Delivery Networks .....	8
3.5 Compute-Fabric-as-a-Service (CFaaS) .....	8
4 Requirements Derived from Use Cases .....	9
5 Enabling Technologies .....	10
5.1 Application Packaging .....	10
5.2 Service Deployment .....	11
5.3 Compute Inter-Connection at Layer 2 .....	12 NEW text
5.4 Service Routing .....	13
5.5 Constraint-based Forwarding Decisions .....	14 Evolved from service pinning
5.6 Collective Communication .....	14 NEW text (formerly 'opp multicast')
5.7 State Synchronization .....	15
5.8 Dynamic Contracts .....	15
6 Security Considerations .....	15 NEW text (on semantic identifiers)
7 IANA Considerations .....	15
8 Conclusion .....	15

Still included but planned to be moved to new use cases draft (headlines already included there)

Still unclear if this will become ultimately part of joint use case (and requirements) draft but will likely update in next revision with clearer linkage to Section 5

Added more text in a number of sub-sections

## 5.3 Compute Inter-Connection at Layer2

- References to using L2 switching for interconnecting distributed compute resources
  - 5GLAN efforts in 3GPP
  - Edge computing in 5G
- Service routing critical in such distributed (L2) environment, similar to intra-DC service scheduling problem

## 5.5 Constraint-based Forwarding Decisions

- Evolved from previous ‘service pinning’ placeholder in V2 of draft
- Extends Section 5.4 discussion on service routing by including **constraints** into forwarding decision between one or more service instance candidates
- Load/latency may not be the only constraints
  - App/service-specific ones may be needed
- Matching operations in intermediary routers over such constraints may be coordinated across several routers to achieve **service scheduling** capability (across distributed compute resources)
- Section 5.4 already provides references to ongoing work, such as CFN-dyncast, ICNRG work etc. that include constraints in forwarding decision

## 5.6 Collective Communication

- Pattern exhibited in number of micro-service scenarios (outlined in Section 2 – use cases) is not just 1:1 but 1:N, N:1, N:M
- Patterns may be short lived with possibly as short-lived as single requests
- Solutions required for supporting such spontaneous formation of multipoint relations
- References to be added to ongoing work in this space (e.g., BIER)

# Future Plans

- Move use cases into updated use case draft
  - > Use clear references in various remaining sub-sections
- Link requirements more clearly to Section 5 sub-sections
- Fill in missing sub-sections for Section 5
- More clearly link to other COIN drafts in relevant areas, e.g., computing frameworks, programmable forwarding nodes
- **Adopt as RG draft?**