



Innovative R&D by NTT

# P4/Stratum活用に向けた プログラマブルスイッチの実装・検討

2018年10月19日

日本電信電話株式会社

NTTネットワークサービスシステム研究所

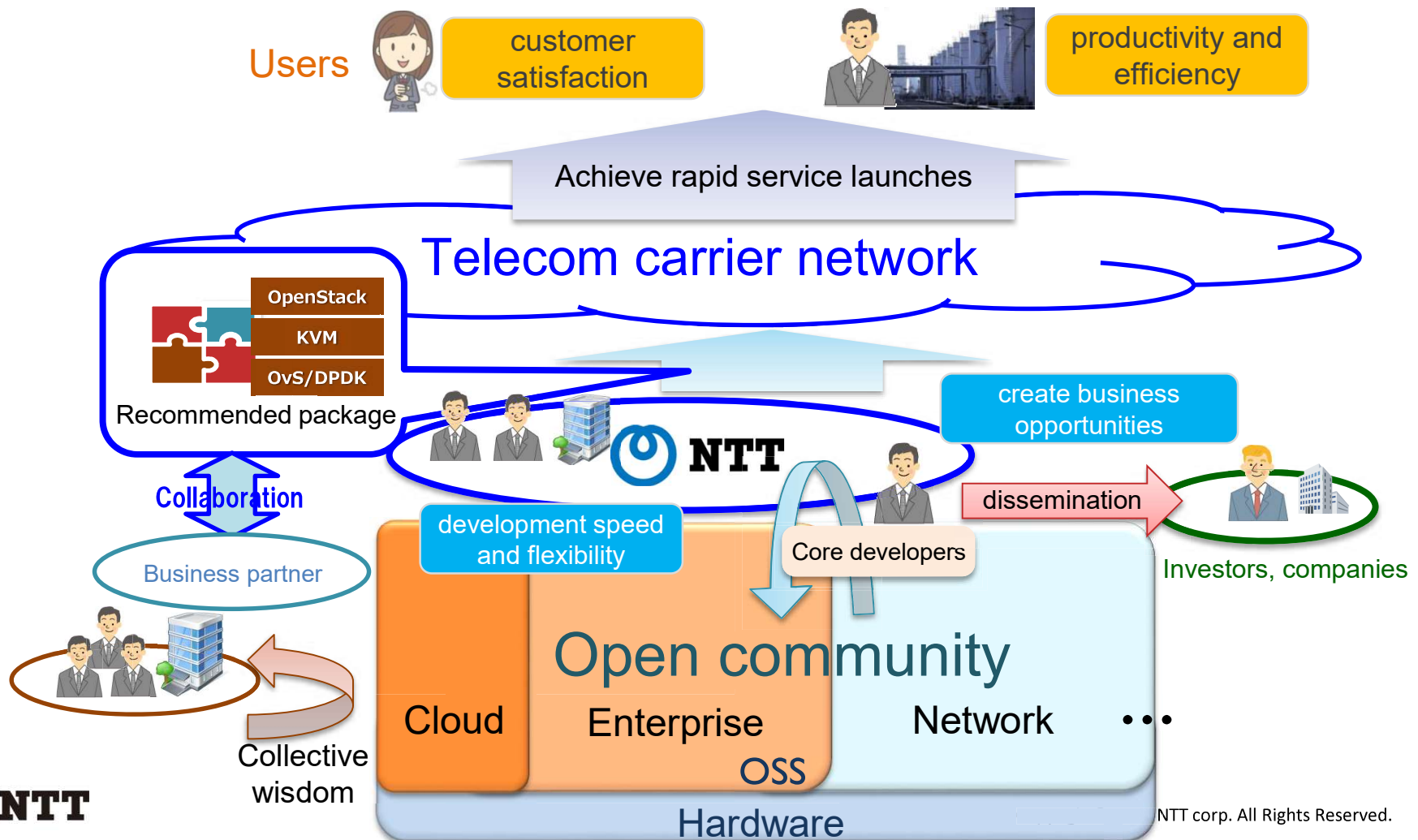
武井 勇樹

## 本日のアジェンダ

- 背景
- P4/Stratumの概要
- NTTのユースケース
- P4スイッチ/P4Runtimeの検討状況
- 今後の取り組み
- まとめ

# NTTネットワークのR&Dビジョン

NTTは、「プロダクトのオープン化」と「OSSへの要件インプット」の両面からコミュニティ活動を実施し、通信キャリアネットワークへの適用を検討している。



# ONF Reference Designs

- NTTではONF (Open Networking Foundation)に参加し、議論を実施している。
- ONFの新たな活動の方向性として、以下の技術分野に中心に取り組む方針を発表(2018.6)。我々はRD-UPANにおいて、P4/Stratumの通信キャリアにおけるユースケースを議論している。

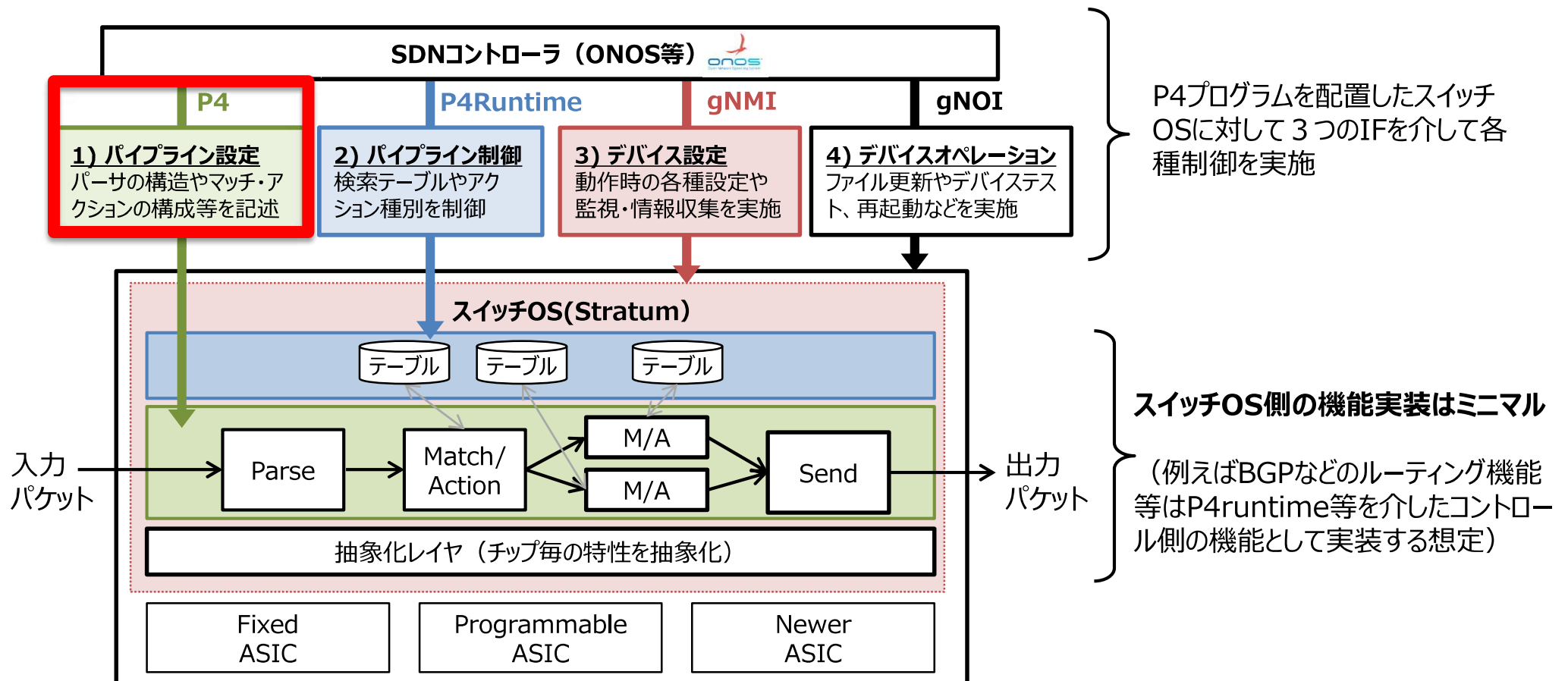
	概要	ベースとなるONFの取り組み	オペレータ
<b>RD-SEBA</b> (SDN Enabled Broadband Assess)	Lightweight reference design supporting a multitude of virtualized access technologies at the edge of the carrier network, including PON, G.Fast, DOCSYS and more.	R-CORD	AT&T, Deutsche Telekom, <b>NTT</b> , Turk Telekom
<b>RD-Trellis</b>	SDN-native spine-leaf data center fabric optimized for edge applications.	Trellis	Comcast
<b>RD-UPAN</b> (Unified, Programmable & Automated Network)	Next generation SDN reference design, leveraging P4 to enable flexible data plane programmability and network embedded VNF acceleration	Stratum	China Unicom, Deutsche Telekom, Google, <b>NTT</b> , Turk Telekom 我々の取り組み領域
<b>RD-ODTN</b> (Open Disaggregated Transport Network)	Open multi-vendor optical networks	ODTN	China Unicom, Comcast, <b>NTT</b>

NTTグループからはRD-SEBA(vOLT)、RD-UPAN(エッジ/P4/Stratum)、RD-ODTN(ODTN)に参加

# P4/Stratumの概要

# Stratumとは

- Stratumは、SDN集中制御を前提としたミニマルな機能実装を指向したスイッチOSである。
- 異なるベンダチップを搭載したスイッチであっても統一的手段で制御するためにP4Runtime・gNMI・gNOIと呼ぶインタフェースを規定。

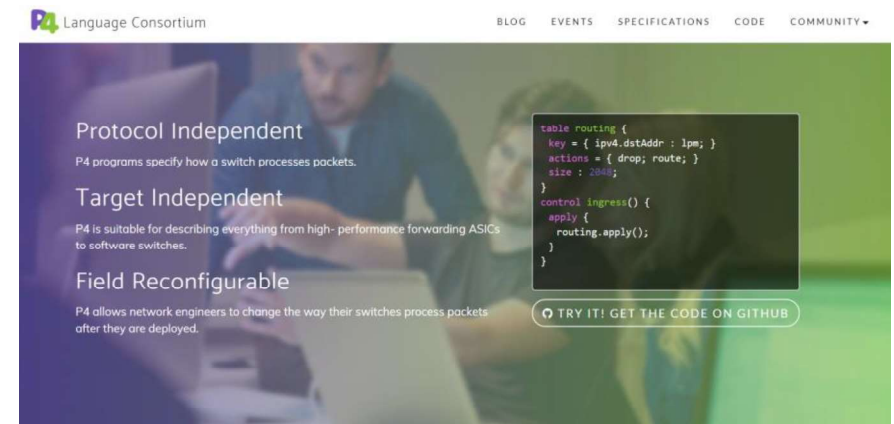


# P4とは



NTT Confidential

- P4 :  
High-level Language for  
Programming Protocol-independent Packet Processors の略
- パケットを受け取ってから、その中身に応じて処理を行い送り出すまで(データプレーン)の動作を記述するためのプログラミング言語
- 『特定のプロトコルに依存せず、  
パケット処理の独自仕様を実現できる』  
(Protocol-independent)
- 『多様なハードウェアの仕様を意識せず、  
パケットの処理の内容を実装できる』  
(Target Independent)



<https://p4.org/>より引用

# P4対応デバイス

- CPU、FPGA、NPU、スイッチASICといった多様なデバイス種別に対応している。
- P4には現在 2 つの仕様バージョン(P4<sub>14</sub> / P4<sub>16</sub>)が存在し、どちらもメンテナンスが継続。

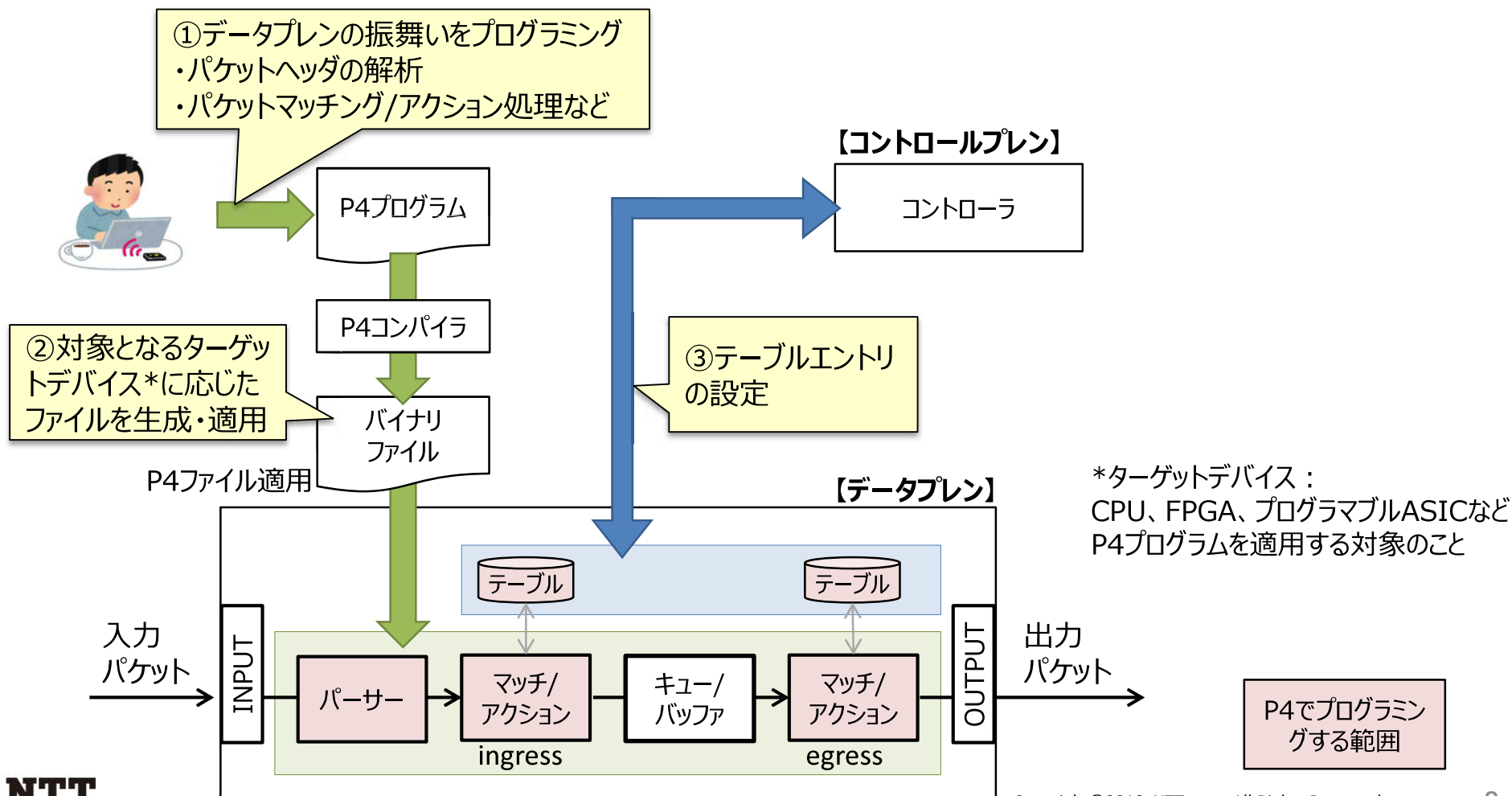
デバイス種別		フレームワーク例	コミュニティ/ベンダ例	サポートP4バージョン
CPU		BMv2	P4.org	P4_14/16
Smart NIC	FPGA	SDNet	Xilinx	P4_16
		NP4	Netcope	P4_14
	NPU	Agilio	Netronome	P4_14/16
Switch-ASIC		Tofino	Barefoot	P4_14/16

※2018.10時点



# P4でパケット転送するまでの基本的な流れ

- データプレンの振舞いを定義するプログラム（P4プログラム）を準備し、コンパイル（①②）
- パケット転送に必要なテーブルエントリを設定し、パケットを転送（③）

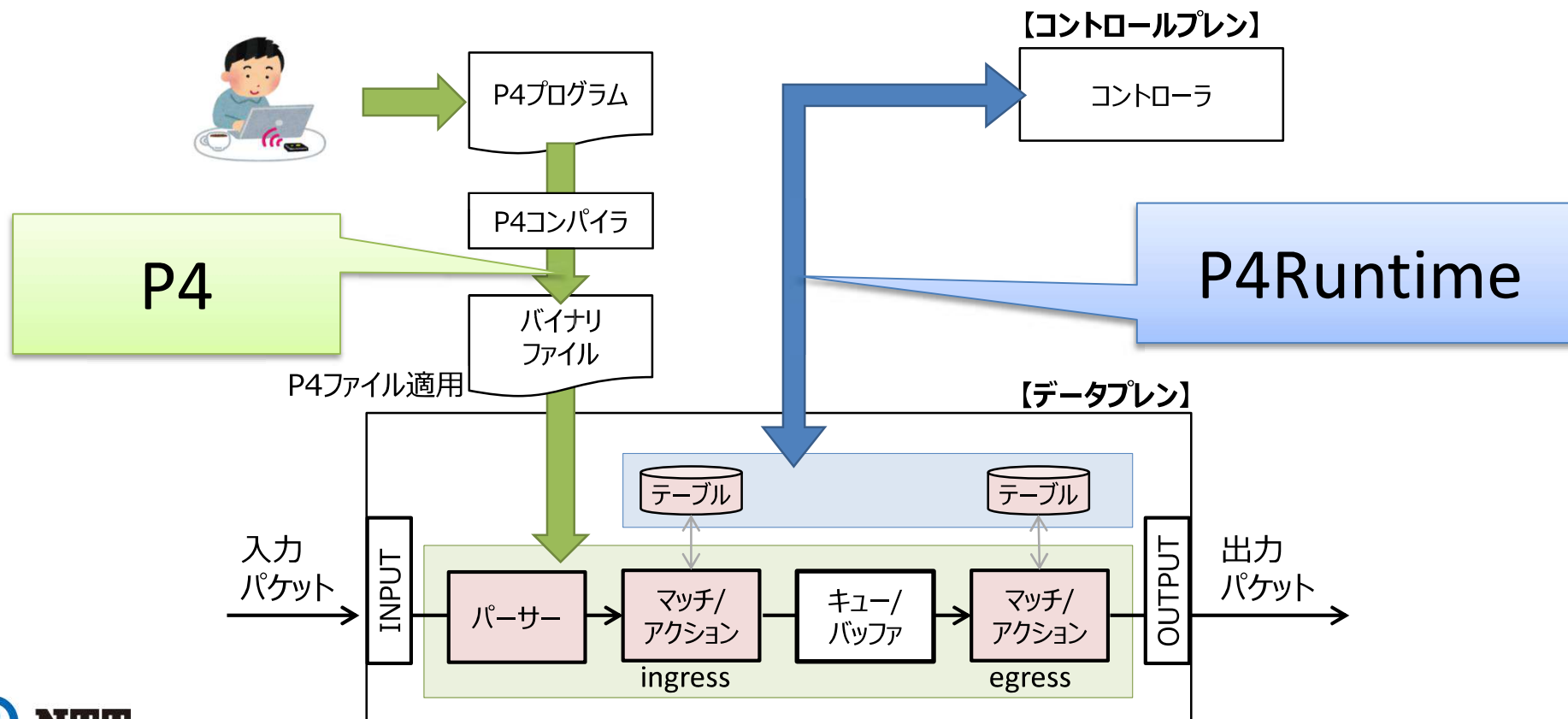


# P4とP4 runtime

P4 : デバイスを抽象化し、スイッチがどのようにパケットを処理するか定義。

パイプラインの指定(どのフィールドが一致するか、どんなアクションを実行するか等)。

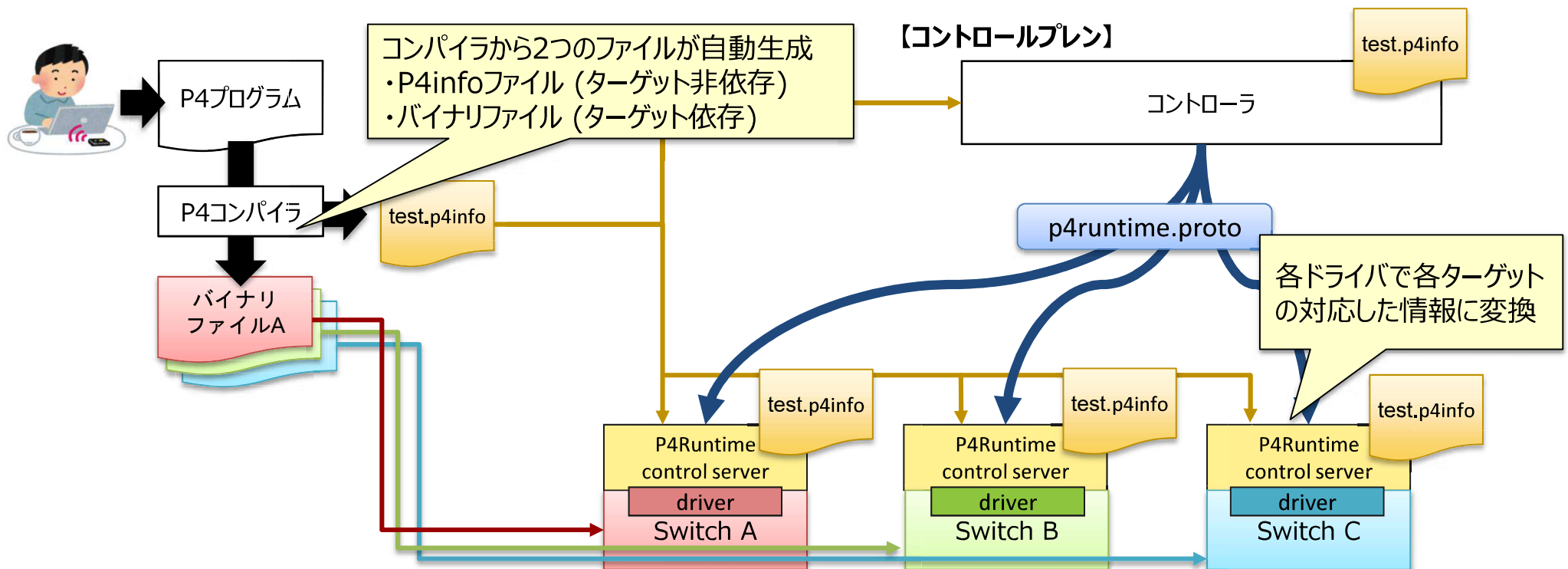
P4Runtime : 固定的な機能のスイッチであるか、プログラマブルであるかに関係なく、P4言語で動作が指定されている機能を制御するために使用するAPI。



# P4処理の流れについて

コンパイラから2つのファイルが自動生成。

- ① P4infoファイル (ターゲット非依存) : コントロールプレーンから制御するために必要なテーブルIDなどの対応関係が格納。外部コントローラはこのファイルを参照しながらターゲットデバイスを制御。
- ② バイナリファイル (ターゲット依存) : 対象となるターゲットデバイスに応じたデータプレーンファイル。



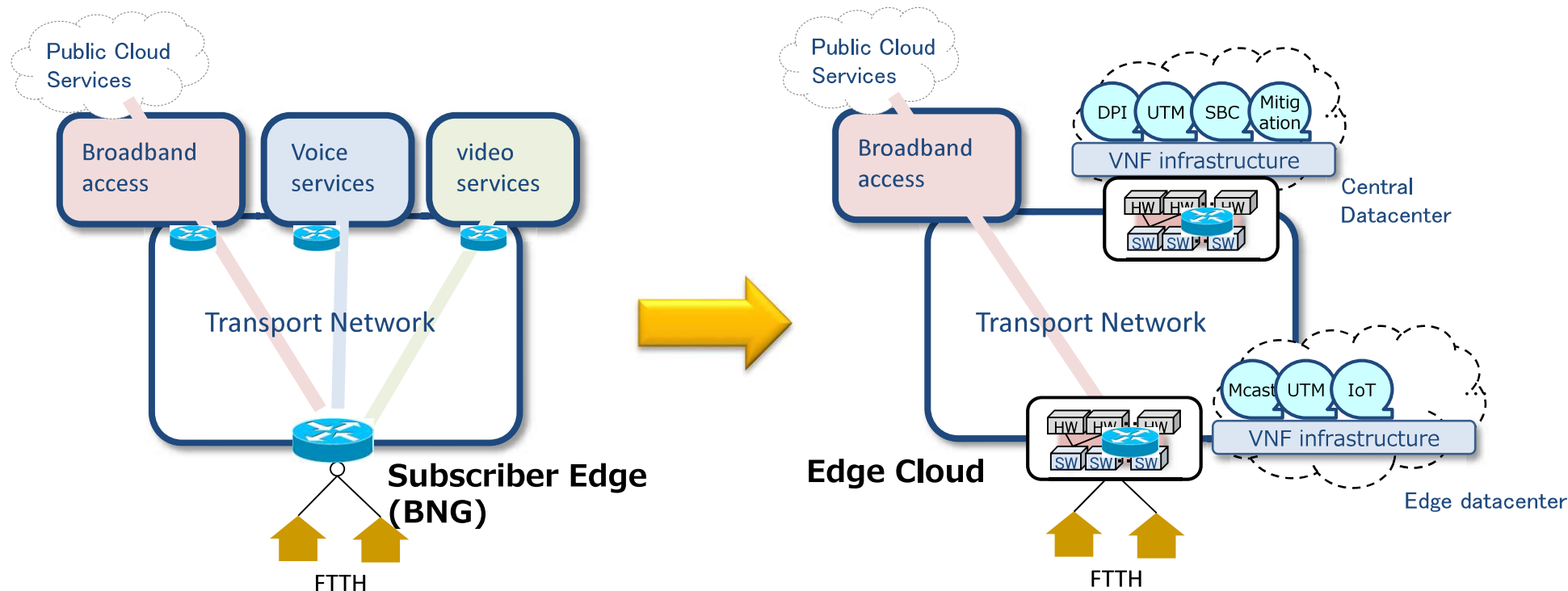


NTT Confidential

# NTTのユースケース

# 将来ネットワークのイメージ

- 将来のネットワークにおいては、サービスの変更・拡張に対して、経済的かつ柔軟に対応するため、プログラマビリティが要求されると考えられる。
- この実現に向けて、P4/Stratumを活用できないか検討している。



- トリプルプレイサービス
- 加入者向けに最適化

- サービスの変更・拡張にあわせ、ネットワーク機能を迅速に提供
- 新しい事業を始めたい事業者や企業に応じたカスタマイズが可能

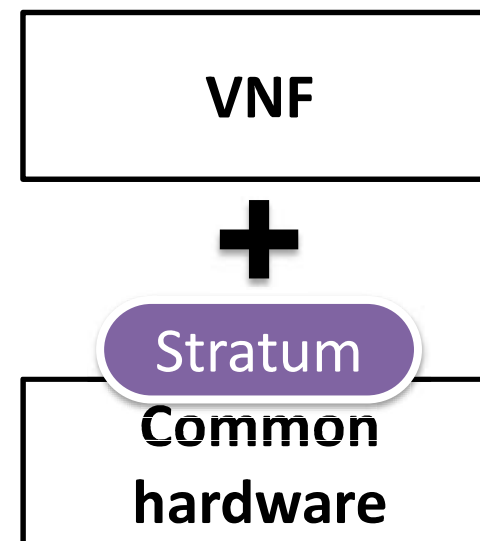
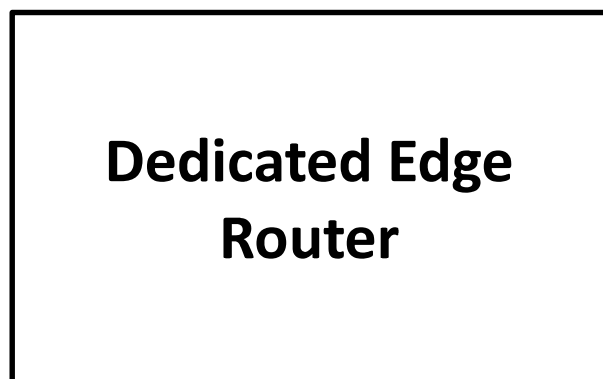
# P4/Stratumのユースケース

- 我々はP4/Stratumを活用のユースケースとして、2つのアプローチを検討。
  - ① BNGのディスアグリゲーション (VNFオフロード)
  - ② 通信品質の可視化 (In-band Network Telemetry)

効用	BNGのディスアグリゲーション (VNFのオフロード)	通信品質の可視化 (In-band Network Telemetry)
効用	NWの経済化	新サービス提供
解決する課題	現状の仮想化NWでは、高度な性能を出すためにはサーバ増設が必要あり、CAPEXが増加してしまう	現状のNW基盤では、低遅延等の新サービスに向けた、厳密な通信品質の管理ができない
実現できること	プログラマブル技術により、サーバ上の機能の一部をスイッチへオフロード可能となり、サーバの増設台数を抑えることができる	プログラマブル技術により、細かい単位で情報を取得可能

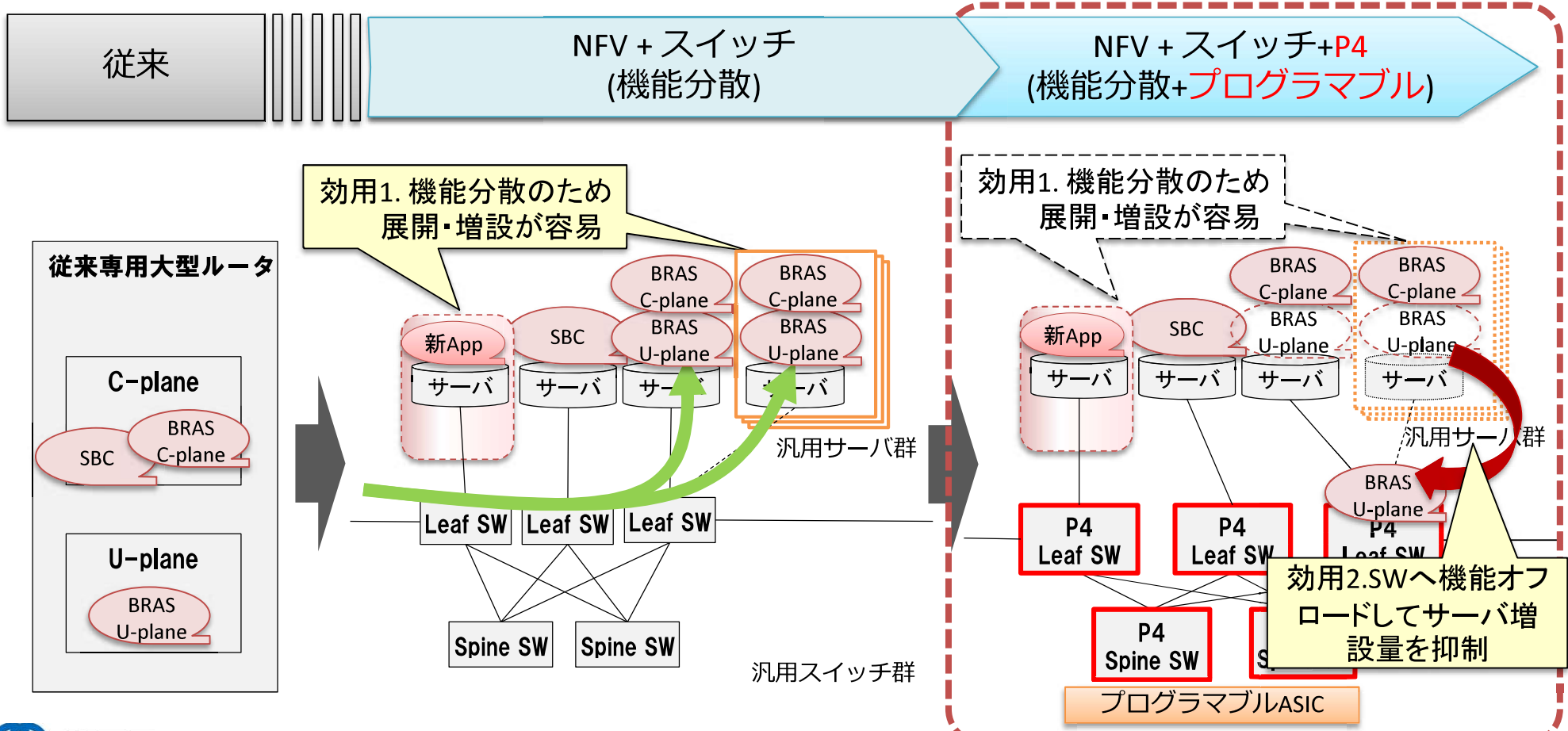
## ユースケース1 BNGのディスアグリゲーション(1/2)

- ネットワークの柔軟性・スケール性の向上のため、既存の大型ルータを汎用装置とVNF技術による実現を目指している。



# ユースケース 1 BNGのディスアグリゲーション(2/2)

汎用サーバだけではなくホワイトボックス側にも自らが定義した機能を自在に分散配備できる。  
 効用 1. 必要な設備を必要なだけ増設すれば良いので、新たなサービスを迅速かつ安価に展開可能。  
 効用 2. 汎用サーバ上で仮想化している機能の内、ボトルネックだけをスイッチ側へ移行可能。

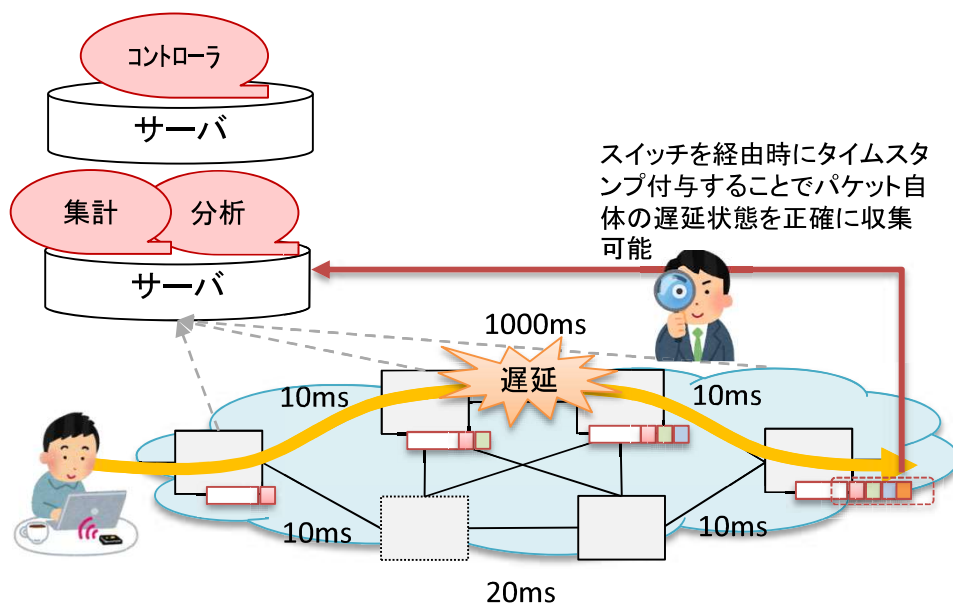




## ユースケース2 通信品質の可視化

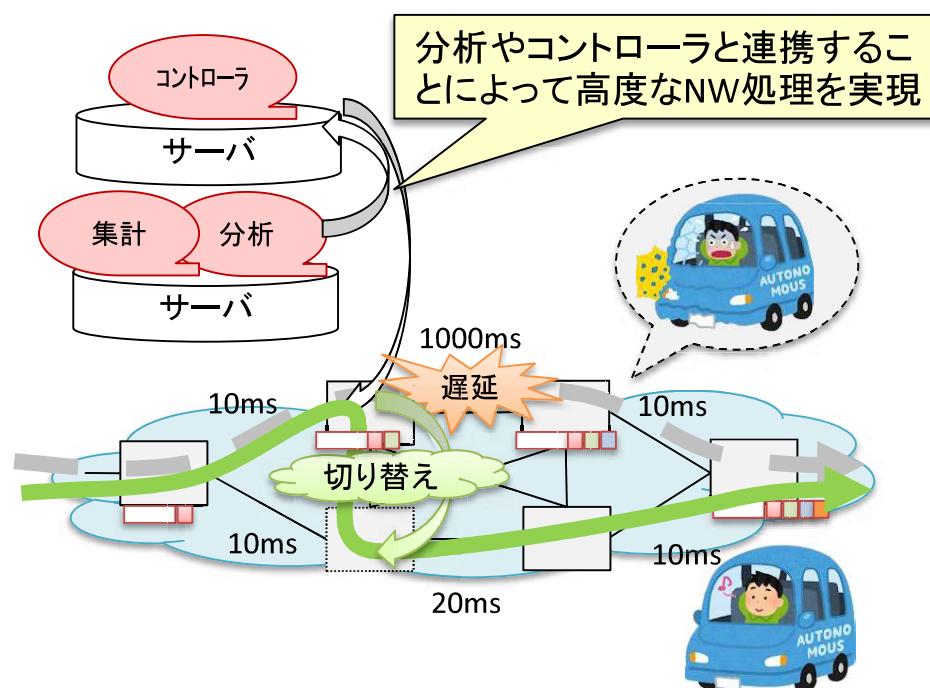
- 情報取得を装置単位からパケット単位にすることができ、通信品質を満たしているかを明確に可視化。
- 細かい単位で情報を取得できるので、精緻な制御を実施可能なインフラを用意できる。

### INTの適応イメージ



遅延情報をパケット単位で可視化  
(低遅延アプリで遅延が発生していないか？等)

### INTの応用例



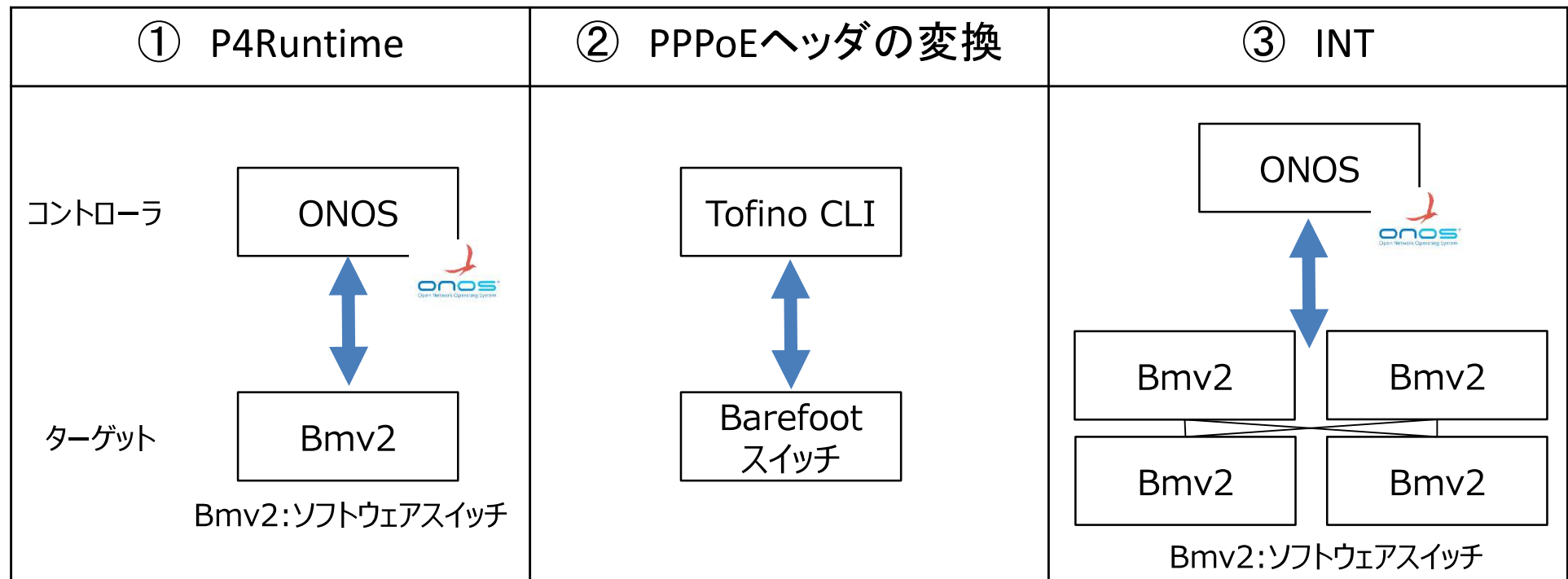
情報を元にパケット単位の緻密な制御を実現  
(遅延発生を予測して経路を迂回する等)

# P4スイッチ/P4Runtimeの検討状況

# 取り組みの全体像

P4/Stratumを用いたユースケース実現に向けて、現在下記の評価・検討を実施。

- ① P4Runtimeの実装状況の確認
- ② P4によるハードウェアオフロードの実現性の確認 (PPPoEヘッダの変換)
- ③ P4によるINTの実現性の確認



# ① P4Runtimeの実装状況 (1/2)

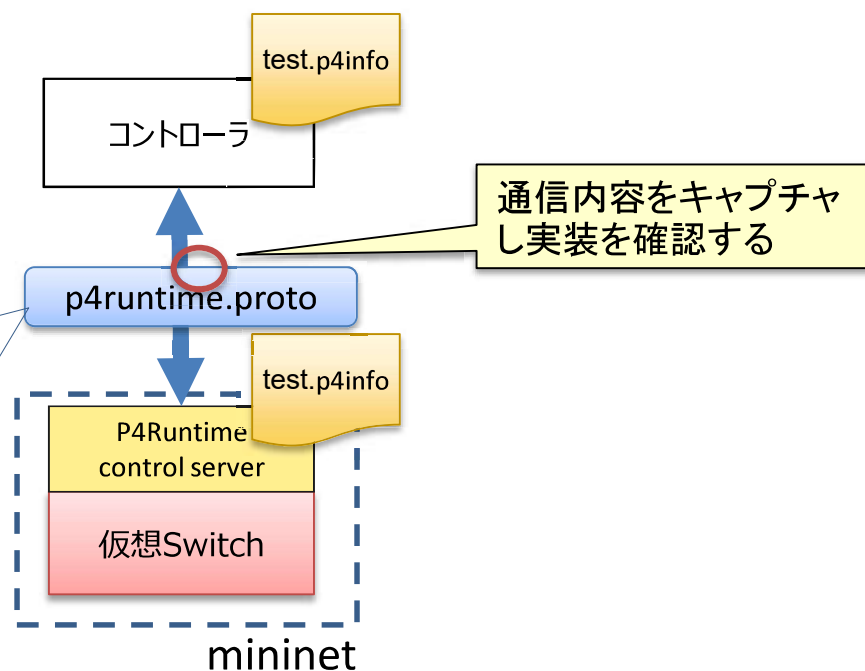
- P4Runtimeは、2018.6に仕様ドラフトがFIXされた。
- ONOSコントローラとデバイス(ソフトウェアスイッチ)間の通信をキャプチャ・Wiresharkで解析することで、P4Runtimeの実装状況を調査する。
- ONOS wikiのチュートリアル(※)を実施。

## p4runtime.protoによる通信内容例

```

service P4Runtime {
  rpc Write(WriteRequest) returns (WriteResponse) {}
  rpc Read(ReadRequest) returns (stream ReadResponse) {}
  rpc SetForwardingPipelineConfig(SetForwardingPipelineConfigRequest)
  returns (SetForwardingPipelineConfigResponse) {}
  rpc GetForwardingPipelineConfig(GetForwardingPipelineConfigRequest)
  returns (GetForwardingPipelineConfigResponse) {}
  rpc StreamChannel(stream StreamMessageRequest)
  returns (stream StreamMessageResponse) {}
}
  
```

<https://github.com/p4lang/p4runtime/blob/master/proto/p4/v1/p4runtime.proto>より



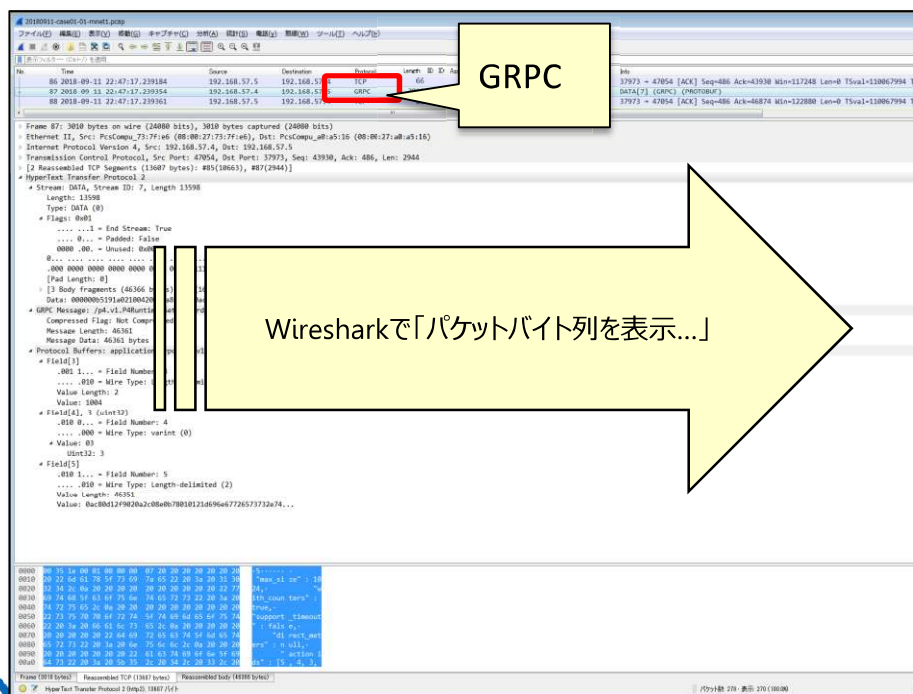
※ONOS wikiのチュートリアルより

<https://wiki.onosproject.org/display/ONOS/P4Runtime+support+in+ONOS>

# ① P4Runtimeの実装状況 (2/2)

- Bmv2環境におけるP4Runtimeの実装状況。
  - ONOS⇔bmv2間でのgRPCベースの通信は実装済みであることを確認。
    - SetForwardingPipelineConfig RequestではBmv2のjsonファイルを送信など。
  - ONOSのGUIからネットワークの構成およびデバイスの情報が表示できる。
- 現在、ONOSとbarefootのスイッチを接続し、同様に通信を解析中。

## 【SetForwardingPipelineConfig Requestの例】



Wiresharkで「パケットバイト列を表示...」

```

1           5. . . . . "max_size" : 1024,
2           "with_counters" : true,
3           "support_timeout" : false,
4           "direct_meters" : null,
5           "action_ids" : [5, 4, 3, 0],
6           "actions" : ["ingress.table0_control.set_egress_port",
"ingress.table0_control.send_to_cpu", "ingress.table0_control.set_next_hop_id",
7           "_drop"],
8           "base_default_next" : "tbl_act_2",
9           "next_tables" : {
10          "tbl_act_2",
11          "ingress.table0_control.send_to_cpu" : "tbl_act_2",
"ingress.table0_control.set_next_hop_id" :
12          "tbl_act_2"
13          }
14          }
15          }
16          }
17          }
18          }
19          }
20          }
21          }
22          }
23          }
24          }
25          }
26          }
27          }
28          }
29          }
30          }
31          }
32          }
33          }
34          }
35          }
36          }
37          }
38          }
39          }
40          }
41          }
42          }
43          }
44          }
45          }
46          }
47          }
48          }
49          }
50          }
51          }
52          }
53          }
54          }
55          }
56          }
57          }
58          }
59          }
60          }
61          }
62          }
63          }
64          }
65          }
66          }
67          }
68          }
69          }
70          }
71          }
72          }
73          }
74          }
75          }
76          }
77          }
78          }
79          }
80          }
81          }
82          }
83          }
84          }
85          }
86          }
87          }
88          }
89          }
90          }
91          }
92          }
93          }
94          }
95          }
96          }
97          }
98          }
99          }
100         }
101         }
102         }
103         }
104         }
105         }
106         }
107         }
108         }
109         }
110         }
111         }
112         }
113         }
114         }
115         }
116         }
117         }
118         }
119         }
120         }
121         }
122         }
123         }
124         }
125         }
126         }
127         }
128         }
129         }
130         }
131         }
132         }
133         }
134         }
135         }
136         }
137         }
138         }
139         }
140         }
141         }
142         }
143         }
144         }
145         }
146         }
147         }
148         }
149         }
150         }
151         }
152         }
153         }
154         }
155         }
156         }
157         }
158         }
159         }
160         }
161         }
162         }
163         }
164         }
165         }
166         }
167         }
168         }
169         }
170         }
171         }
172         }
173         }
174         }
175         }
176         }
177         }
178         }
179         }
180         }
181         }
182         }
183         }
184         }
185         }
186         }
187         }
188         }
189         }
190         }
191         }
192         }
193         }
194         }
195         }
196         }
197         }
198         }
199         }
200         }
201         }
202         }
203         }
204         }
205         }
206         }
207         }
208         }
209         }
210         }
211         }
212         }
213         }
214         }
215         }
216         }
217         }
218         }
219         }
220         }
221         }
222         }
223         }
224         }
225         }
226         }
227         }
228         }
229         }
230         }
231         }
232         }
233         }
234         }
235         }
236         }
237         }
238         }
239         }
240         }
241         }
242         }
243         }
244         }
245         }
246         }
247         }
248         }
249         }
250         }
251         }
252         }
253         }
254         }
255         }
256         }
257         }
258         }
259         }
260         }
261         }
262         }
263         }
264         }
265         }
266         }
267         }
268         }
269         }
270         }
271         }
272         }
273         }
274         }
275         }
276         }
277         }
278         }
279         }
280         }
281         }
282         }
283         }
284         }
285         }
286         }
287         }
288         }
289         }
290         }
291         }
292         }
293         }
294         }
295         }
296         }
297         }
298         }
299         }
300         }
301         }
302         }
303         }
304         }
305         }
306         }
307         }
308         }
309         }
310         }
311         }
312         }
313         }
314         }
315         }
316         }
317         }
318         }
319         }
320         }
321         }
322         }
323         }
324         }
325         }
326         }
327         }
328         }
329         }
330         }
331         }
332         }
333         }
334         }
335         }
336         }
337         }
338         }
339         }
340         }
341         }
342         }
343         }
344         }
345         }
346         }
347         }
348         }
349         }
350         }
351         }
352         }
353         }
354         }
355         }
356         }
357         }
358         }
359         }
360         }
361         }
362         }
363         }
364         }
365         }
366         }
367         }
368         }
369         }
370         }
371         }
372         }
373         }
374         }
375         }
376         }
377         }
378         }
379         }
380         }
381         }
382         }
383         }
384         }
385         }
386         }
387         }
388         }
389         }
390         }
391         }
392         }
393         }
394         }
395         }
396         }
397         }
398         }
399         }
400         }
401         }
402         }
403         }
404         }
405         }
406         }
407         }
408         }
409         }
410         }
411         }
412         }
413         }
414         }
415         }
416         }
417         }
418         }
419         }
420         }
421         }
422         }
423         }
424         }
425         }
426         }
427         }
428         }
429         }
430         }
431         }
432         }
433         }
434         }
435         }
436         }
437         }
438         }
439         }
440         }
441         }
442         }
443         }
444         }
445         }
446         }
447         }
448         }
449         }
450         }
451         }
452         }
453         }
454         }
455         }
456         }
457         }
458         }
459         }
460         }
461         }
462         }
463         }
464         }
465         }
466         }
467         }
468         }
469         }
470         }
471         }
472         }
473         }
474         }
475         }
476         }
477         }
478         }
479         }
480         }
481         }
482         }
483         }
484         }
485         }
486         }
487         }
488         }
489         }
490         }
491         }
492         }
493         }
494         }
495         }
496         }
497         }
498         }
499         }
500         }
501         }
502         }
503         }
504         }
505         }
506         }
507         }
508         }
509         }
510         }
511         }
512         }
513         }
514         }
515         }
516         }
517         }
518         }
519         }
520         }
521         }
522         }
523         }
524         }
525         }
526         }
527         }
528         }
529         }
530         }
531         }
532         }
533         }
534         }
535         }
536         }
537         }
538         }
539         }
540         }
541         }
542         }
543         }
544         }
545         }
546         }
547         }
548         }
549         }
550         }
551         }
552         }
553         }
554         }
555         }
556         }
557         }
558         }
559         }
560         }
561         }
562         }
563         }
564         }
565         }
566         }
567         }
568         }
569         }
570         }
571         }
572         }
573         }
574         }
575         }
576         }
577         }
578         }
579         }
580         }
581         }
582         }
583         }
584         }
585         }
586         }
587         }
588         }
589         }
590         }
591         }
592         }
593         }
594         }
595         }
596         }
597         }
598         }
599         }
600         }
601         }
602         }
603         }
604         }
605         }
606         }
607         }
608         }
609         }
610         }
611         }
612         }
613         }
614         }
615         }
616         }
617         }
618         }
619         }
620         }
621         }
622         }
623         }
624         }
625         }
626         }
627         }
628         }
629         }
630         }
631         }
632         }
633         }
634         }
635         }
636         }
637         }
638         }
639         }
640         }
641         }
642         }
643         }
644         }
645         }
646         }
647         }
648         }
649         }
650         }
651         }
652         }
653         }
654         }
655         }
656         }
657         }
658         }
659         }
660         }
661         }
662         }
663         }
664         }
665         }
666         }
667         }
668         }
669         }
670         }
671         }
672         }
673         }
674         }
675         }
676         }
677         }
678         }
679         }
680         }
681         }
682         }
683         }
684         }
685         }
686         }
687         }
688         }
689         }
690         }
691         }
692         }
693         }
694         }
695         }
696         }
697         }
698         }
699         }
700         }
701         }
702         }
703         }
704         }
705         }
706         }
707         }
708         }
709         }
710         }
711         }
712         }
713         }
714         }
715         }
716         }
717         }
718         }
719         }
720         }
721         }
722         }
723         }
724         }
725         }
726         }
727         }
728         }
729         }
730         }
731         }
732         }
733         }
734         }
735         }
736         }
737         }
738         }
739         }
740         }
741         }
742         }
743         }
744         }
745         }
746         }
747         }
748         }
749         }
750         }
751         }
752         }
753         }
754         }
755         }
756         }
757         }
758         }
759         }
760         }
761         }
762         }
763         }
764         }
765         }
766         }
767         }
768         }
769         }
770         }
771         }
772         }
773         }
774         }
775         }
776         }
777         }
778         }
779         }
780         }
781         }
782         }
783         }
784         }
785         }
786         }
787         }
788         }
789         }
790         }
791         }
792         }
793         }
794         }
795         }
796         }
797         }
798         }
799         }
800         }
801         }
802         }
803         }
804         }
805         }
806         }
807         }
808         }
809         }
810         }
811         }
812         }
813         }
814         }
815         }
816         }
817         }
818         }
819         }
820         }
821         }
822         }
823         }
824         }
825         }
826         }
827         }
828         }
829         }
830         }
831         }
832         }
833         }
834         }
835         }
836         }
837         }
838         }
839         }
840         }
841         }
842         }
843         }
844         }
845         }
846         }
847         }
848         }
849         }
850         }
851         }
852         }
853         }
854         }
855         }
856         }
857         }
858         }
859         }
860         }
861         }
862         }
863         }
864         }
865         }
866         }
867         }
868         }
869         }
870         }
871         }
872         }
873         }
874         }
875         }
876         }
877         }
878         }
879         }
880         }
881         }
882         }
883         }
884         }
885         }
886         }
887         }
888         }
889         }
890         }
891         }
892         }
893         }
894         }
895         }
896         }
897         }
898         }
899         }
900         }
901         }
902         }
903         }
904         }
905         }
906         }
907         }
908         }
909         }
910         }
911         }
912         }
913         }
914         }
915         }
916         }
917         }
918         }
919         }
920         }
921         }
922         }
923         }
924         }
925         }
926         }
927         }
928         }
929         }
930         }
931         }
932         }
933         }
934         }
935         }
936         }
937         }
938         }
939         }
940         }
941         }
942         }
943         }
944         }
945         }
946         }
947         }
948         }
949         }
950         }
951         }
952         }
953         }
954         }
955         }
956         }
957         }
958         }
959         }
960         }
961         }
962         }
963         }
964         }
965         }
966         }
967         }
968         }
969         }
970         }
971         }
972         }
973         }
974         }
975         }
976         }
977         }
978         }
979         }
980         }
981         }
982         }
983         }
984         }
985         }
986         }
987         }
988         }
989         }
990         }
991         }
992         }
993         }
994         }
995         }
996         }
997         }
998         }
999         }
1000        }
1001        }
1002        }
1003        }
1004        }
1005        }
1006        }
1007        }
1008        }
1009        }
1010        }
1011        }
1012        }
1013        }
1014        }
1015        }
1016        }
1017        }
1018        }
1019        }
1020        }
1021        }
1022        }
1023        }
1024        }
1025        }
1026        }
1027        }
1028        }
1029        }
1030        }
1031        }
1032        }
1033        }
1034        }
1035        }
1036        }
1037        }
1038        }
1039        }
1040        }
1041        }
1042        }
1043        }
1044        }
1045        }
1046        }
1047        }
1048        }
1049        }
1050        }
1051        }
1052        }
1053        }
1054        }
1055        }
1056        }
1057        }
1058        }
1059        }
1060        }
1061        }
1062        }
1063        }
1064        }
1065        }
1066        }
1067        }
1068        }
1069        }
1070        }
1071        }
1072        }
1073        }
1074        }
1075        }
1076        }
1077        }
1078        }
1079        }
1080        }
1081        }
1082        }
1083        }
1084        }
1085        }
1086        }
1087        }
1088        }
1089        }
1090        }
1091        }
1092        }
1093        }
1094        }
1095        }
1096        }
1097        }
1098        }
1099        }
1100        }
1101        }
1102        }
1103        }
1104        }
1105        }
1106        }
1107        }
1108        }
1109        }
1110        }
1111        }
1112        }
1113        }
1114        }
1115        }
1116        }
1117        }
1118        }
1119        }
1120        }
1121        }
1122        }
1123        }
1124        }
1125        }
1126        }
1127        }
1128        }
1129        }
1130        }
1131        }
1132        }
1133        }
1134        }
1135        }
1136        }
1137        }
1138        }
1139        }
1140        }
1141        }
1142        }
1143        }
1144        }
1145        }
1146        }
1147        }
1148        }
1149        }
1150        }
1151        }
1152        }
1153        }
1154        }
1155        }
1156        }
1157        }
1158        }
1159        }
1160        }
1161        }
1162        }
1163        }
1164        }
1165        }
1166        }
1167        }
1168        }
1169        }
1170        }
1171        }
1172        }
1173        }
1174        }
1175        }
1176        }
1177        }
1178        }
1179        }
1180        }
1181        }
1182        }
1183        }
1184        }
1185        }
1186        }
1187        }
1188        }
1189        }
1190        }
1191        }
1192        }
1193        }
1194        }
1195        }
1196        }
1197        }
1198        }
1199        }
1200        }
1201        }
1202        }
1203        }
1204        }
1205        }
1206        }
1207        }
1208        }
1209        }
1210        }
1211        }
1212        }
1213        }
1214        }
1215        }
1216        }
1217        }
1218        }
1219        }
1220        }
1221        }
1222        }
1223        }
1224        }
1225        }
1226        }
1227        }
1228        }
1229        }
1230        }
1231        }
1232        }
1233        }
1234        }
1235        }
1236        }
1237        }
1238        }
1239        }
1240        }
1241        }
1242        }
1243        }
1244        }
1245        }
1246        }
1247        }
1248        }
1249        }
1250        }
1251        }
1252        }
1253        }
1254        }
1255        }
1256        }
1257        }
1258        }
1259        }
1260        }
1261        }
1262        }
1263        }
1264        }
1265        }
1266        }
1267        }
1268        }
1269        }
1270        }
1271        }
1272        }
1273        }
1274        }
1275        }
1276        }
1277        }
1278        }
1279        }
1280        }
1281        }
1282        }
1283        }
1284        }
1285        }
1286        }
1287        }
1288        }
1289        }
1290        }
1291        }
1292        }
1293        }
1294        }
1295        }
1296        }
1297        }
1298        }
1299        }
1300        }
1301        }
1302        }
1303        }
1304        }
1305        }
1306        }
1307        }
1308        }
1309        }
1310        }
1311        }
1312        }
1313        }
1314        }
1315        }
1316        }
1317        }
1318        }
1319        }
1320        }
1321        }
1322        }
1323        }
1324        }
1325        }
1326        }
1327        }
1328        }
1329        }
1330        }
1331        }
1332        }
1333        }
1334        }
1335        }
1336        }
1337        }
1338        }
1339        }
1340        }
1341        }
1342        }
1343        }
1344        }
1345        }
1346        }
1347        }
1348        }
1349        }
1350        }
1351        }
1352        }
1353        }
1354        }
1355        }
1356        }
1357        }
1358        }
1359        }
1360        }
1361        }
1362        }
1363        }
1364        }
1365        }
1366        }
1367        }
1368        }
1369        }
1370        }
1371        }
1372        }
1373        }
1374        }
1375        }
1376        }
1377        }
1378        }
1379        }
1380        }
1381        }
1382        }
1383        }
1384        }
1385        }
1386        }
1387        }
1388        }
1389        }
1390        }
1391        }
1392        }
1393        }
1394        }
1395        }
1396        }
1397        }
1398        }
1399        }
1400        }
1401        }
1402        }
1403        }
1404        }
1405        }
1406        }
1407        }
1408        }
1409        }
1410        }
1411        }
1412        }
1413        }
1414        }
1415        }
1416        }
1417        }
1418        }
1419        }
1420        }
1421        }
1422        }
1423        }
1424        }
1425        }
1426        }
1427        }
1428        }
1429        }
1430        }
1431        }
1432        }
1433        }
1434        }
1435        }
1436        }
1437        }
1438        }
1439        }
1440        }
1441        }
1442        }
1443        }
1444        }
1445        }
1446        }
1447        }
1448        }
1449        }
1450        }
1451        }
1452        }
1453        }
1454        }
1455        }
1456        }
1457        }
1458        }
1459        }
1460        }
1461        }
1462        }
1463        }
1464        }
1465        }
1466        }
1467        }
1468        }
1469        }
1470        }
1471        }
1472        }
1473        }
1474        }
1475        }
1476        }
1477        }
1478        }
1479        }
1480        }
1481        }
1482        }
1483        }
1484        }
1485        }
1486        }
1487        }
1488        }
1489        }
1490        }
1491        }
1492        }
1493        }
1494        }
1495        }
1496        }
1497        }
1498        }
1499        }
1500        }
1501        }
1502        }
1503        }
1504        }
1505        }
1506        }
1507        }
1508        }
1509        }
1510        }
1511        }
1512        }
1513        }
1514        }
1515        }
1516        }
1517        }
1518        }
1519        }
1520        }
1521        }
1522        }
1523        }
1524        }
1525        }
1526        }
1527        }
1528        }
1529        }
1530        }
1531        }
1532        }
1533        }
1534        }
1535        }
1536        }
1537        }
1538        }
1539        }
1540        }
1541        }
1542        }
1543        }
1544        }
1545        }
1546        }
1547        }
1548        }
1549        }
1550        }
1551        }
1552        }
1553        }
1554        }
1555        }
1556        }
1557        }
1558        }
1559        }
1560        }
1561        }
1562        }
1563        }
1564        }
1565        }
1566        }
1567        }
1568        }
1569        }
1570        }
1571        }
1572        }
1573        }
1574        }
1575        }
1576        }
1577        }
1578        }
1579        }
1580        }
1581        }
1582        }
1583        }
1584        }
1585        }
1586        }
1587        }
1588        }
1589        }
1590        }
1591        }
1592        }
1593        }
1594        }
1595        }
1596        }
1597        }
1598        }
1599        }
1600        }
1601        }
1602        }
1603        }
1604        }
1605        }
1606        }
1607        }
1608        }
1609        }
1610        }
1611        }
1612        }
1613        }
1614        }
1615        }
1616        }
1617        }
1618        }
1619        }
1620        }
1621        }
1622        }
1623        }
1624        }
1625        }
1626        }
1627        }
1628        }
1629        }
1630        }
1631        }
1632        }
1633        }
1634        }
1635        }
1636        }
1637        }
1638        }
1639        }
1640        }
1641        }
1642        }
1643        }
1644        }
1645        }
1646        }
1647        }
1648        }
1649        }
1650        }
1651        }
1652        }
1653        }
1654        }
1655        }
1656        }
1657        }
1658        }
1659        }
1660        }
1661        }
1662        }
1663        }
1664        }
1665        }
1666        }
1667        }
1668        }
1669        }
1670        }
1671        }
1672        }
1673        }
1674        }
1675        }
1676        }
1677        }
1678        }
1679        }
1680        }
1681        }
1682        }
1683        }
1684        }
1685        }
1686        }
1687        }
1688        }
1689        }
1690        }
1691        }
1692        }
1693        }
1694        }
1695        }
1696        }
1697        }
1698        }
1699        }
1700        }
1701        }
1702        }
1703        }
1704        }
1705        }
1706        }
1707        }
1708        }
1709        }
1710        }
1711        }
1712        }
1713        }
1714        }
1715        }
1716        }
1717        }
1718        }
1719        }
1720        }
1721        }
1722        }
1723        }
1724        }
1725        }
1726        }
1727        }
1728        }
1729        }
1730        }
1731        }
1732        }
1733        }
1734        }
1735        }
1736        }
1737        }
1738        }
1739        }
1740        }
1741        }
1742        }
1743        }
1744        }
1745        }
1746        }
1747        }
1748        }
1749        }
1750        }
1751        }
1752        }
1753        }
1754        }
1755        }
1756        }
1757        }
1758        }
1759        }
1760        }
1761        }
1762        }
1763        }
1764        }
1765        }
1766        }
1767        }
1768        }
1769        }
1770        }
1771        }
1772        }
1773        }
1774        }
1775        }
1776        }
1777        }
1778        }
1779        }
1780        }
1781        }
1782        }
1783        }
1784        }
1785        }
1786        }
1787        }
1788        }
1789        }
1790        }
1791        }
1792        }
1793        }
1794        }
1795        }
1796        }
1797        }
1798        }
1799        }
1800        }
1801        }
1802        }
1803        }
1804        }
1805        }
1806        }
1807        }
1808        }
1809        }
1810        }
1811        }
1812        }
1813        }
1814        }
1815        }
1816        }
1817        }
1818        }
1819        }
1820        }
1821        }
1822        }
1823        }
1824        }
1825        }
1826        }
1827        }
1828        }
1829        }
1830        }
1831        }
1832        }
1833        }
1834        }
1
```

## ② P4によるPPPoEヘッダの変換の実装 (1/2)

- PPPoEパケットをL2TPパケットに変換するオリジナルコードをP4\_14で作成。  
→ 非常に簡単に実装が可能 (400行程度)
- Barefoot SDEを用いてスイッチ上に実装。

```

PPPoE.p4 (P4_14)
/***** headers *****/
...
header_type pppoe_t {
  fields {
    ver      : 4;
    pktType  : 4;
    code     : 8;
    sessionId : 16;
    payloadLength : 16;
  }
}
...
/***** parser *****/
...
/***** action *****/
...
/***** table *****/
...
/***** control *
...
  
```

400行程度  
(改行・コメント含)

PPPoE.p4

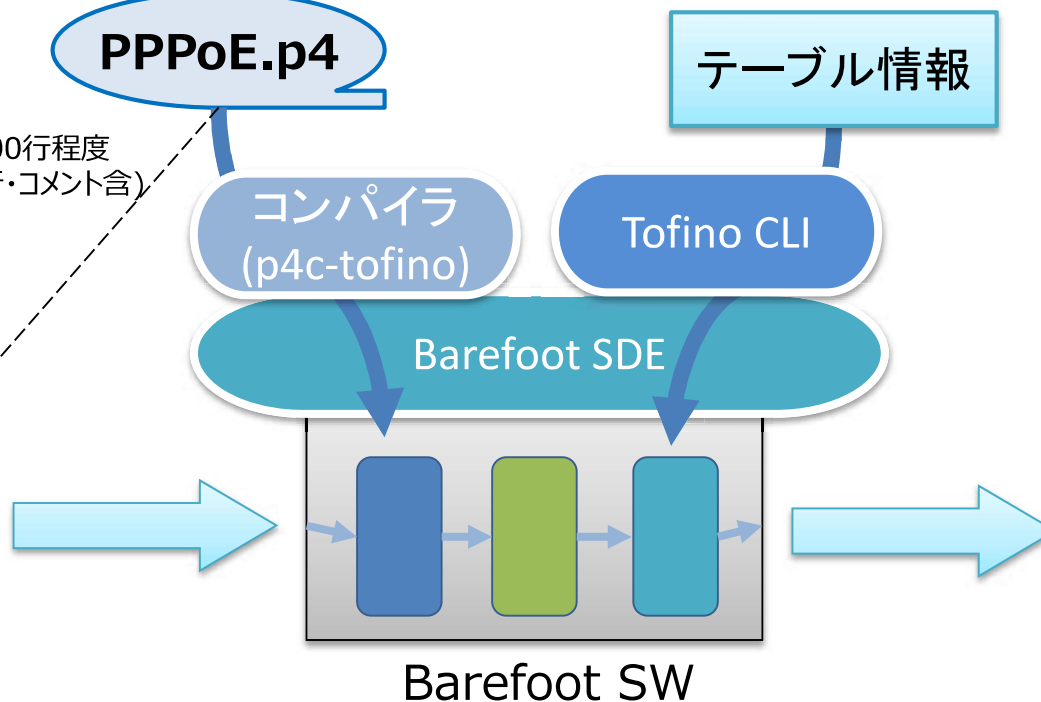
コンパイラ  
(p4c-tofino)

Tofino CLI

テーブル情報

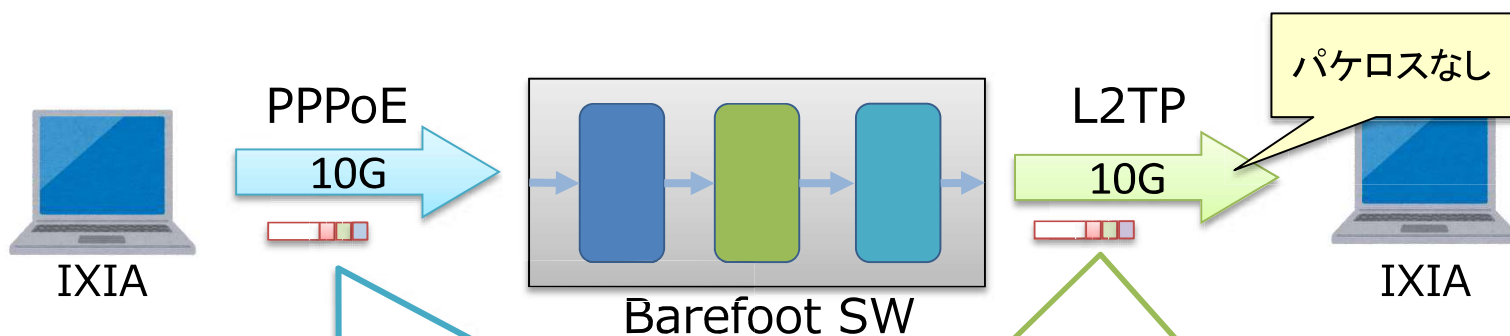
Barefoot SDE

Barefoot SW



## ② P4によるPPPoEヘッダの変換の実装 (2/2)

- 試験機より疑似パケットを印加し、実装した機能が動作するか確認。
- 下記のように、ハードウェアスイッチ上でパケットの変換動作が可能。



```

> Frame 1: 974 bytes on wire (7792 bits), 974 bytes captured (7792 bits) on interface 0
> Ethernet II, Src: MS-NLB-PhysServer-06_0a:0e:ff:f3 (02:06:0a:0e:ff:f3), Dst: Rivertre_0
> 802.1Q Virtual LAN, PRI: 0, DEI: 0, ID: 1001
> PPP-over-Ethernet Session
  0001 .... = Version: 1
  .... 0001 = Type: 1
  Code: Session Data (0x00)
  Session ID: 0x0001
  Payload Length: 946
> Point-to-Point Protocol
> Internet Protocol Version 4, Src: 2.2.2.2, Dst: 1.1.1.1
> User Datagram Protocol, Src Port: 60002, Dst Port: 60001
> Data (916 bytes)
  
```

PPPoE

```

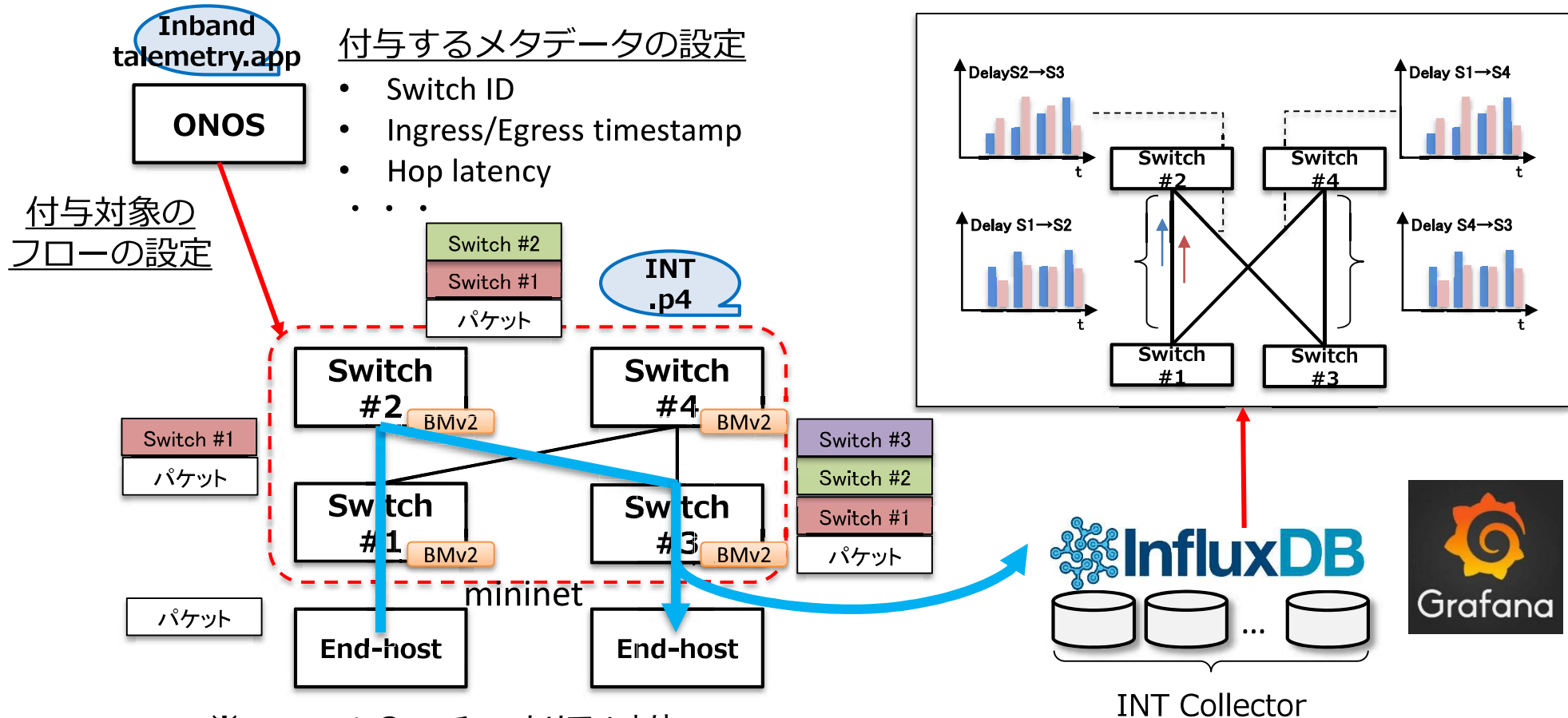
> Frame 1: 1000 bytes on wire (8000 bits), 1000 bytes captured (8000 bits) on interface 0
> Ethernet II, Src: Cisco_00:00:01 (00:12:01:00:00:01), Dst: AlaxalaN_e3:0b:bb (00:12:e2:e
> Internet Protocol Version 4, Src: 10.3.0.2, Dst: 100.0.0.1
> User Datagram Protocol, Src Port: 1701, Dst Port: 1701
> Layer 2 Tunneling Protocol
  > Packet Type: Data Message Tunnel Id=61847 Session Id=64798
  Tunnel ID: 61847
  Session ID: 64798
> Point-to-Point Protocol
> Internet Protocol Version 4, Src: 2.2.2.2, Dst: 1.1.1.1
> User Datagram Protocol, Src Port: 60002, Dst Port: 60001
> Data (916 bytes)
  
```

L2TP



### ③ P4によるINTの実現性の確認 (1/3)

- INT(In-band Network Telemetry) とONOSを用いたINT実装を確認。
- ONOS wikiのINTチュートリアルを使用。

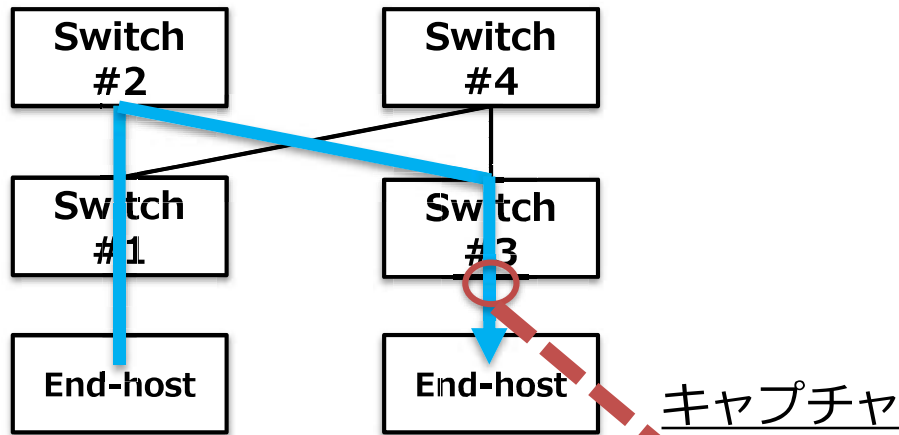


※ ONOS wikiのINTチュートリアルより

<https://wiki.onosproject.org/display/ONOS/%28INT%29+In+Band+Network+Telemetry+with+ONOS+and+P4>



### ③ P4によるINTの実現性の確認 (2/3)



```

0000 ff ff ff ff ff ff a4 23 05 00 00 00 08 00 45 00 .....#.....E.
0010 01 b0 00 00 00 00 ff 11 00 00 0a c0 13 b4 7f 00 .....#.....E.
0020 00 01 00 00 d4 31 01 9c 00 00 00 20 00 01 00 00 .....1.....
0030 00 00 3c 6c 84 69 00 00 00 00 02 05 00 00 00 00 ...<1.i.....
0040 02 26 08 00 45 04 01 7a 29 07 40 00 3f 11 f8 65 ..&..E..z )-@..?..e
0050 0a 00 02 01 0a 00 03 02 d0 4e 13 89 01 66 1a 16 .....N...f..
0060 01 00 19 00 00 07 40 03 fd 00 00 00 00 00 00 cd .....@.....
0070 00 00 00 05 00 00 00 0a 00 00 00 00 3c 6c 84 69 .....<1.i.....
0080 3c 6c 84 73 00 00 00 00 00 00 00 e2 00 01 00 02 <1.s.....
0090 00 00 00 19 00 00 00 00 3c 63 f7 c7 3c 63 .....
00a0 00 00 00 00 00 00 00 cc 00 03 00 01 00 00 .....
00b0 00 00 00 00 3c 73 41 25 3c 73 41 2c 00 00 .....
00c0 11 13 89 00 00 29 e9 c8 5b c4 11 a9 00 0c .....
00d0 00 00 00 00 00 00 00 01 00 00 13 89 00 00 .....
00e0 00 10 00 00 ff f0 bd c0 36 37 38 39 30 31 .....
00f0 34 35 36 37 38 39 30 31 32 33 34 35 36 37 .....
0100 30 31 32 33 34 35 36 37 38 39 30 31 32 33 .....
0110 36 37 38 39 30 31 32 33 34 35 36 37 38 39 .....
0120 32 33 34 35 36 37 38 39 30 31 32 33 34 35 .....
  
```

Eth + IP + UDP ヘッダ

INT メタデータ

- Switch #3 INT META DATA
- Switch #2 INT META DATA
- Switch #1 INT META DATA

ペイロード

```

Frame 1: 488 bytes on wire (3904 bits), 488 bytes captured (3904 bits)
Ethernet II, Src: OpenNetw_00:00:00 (a4:23:05:00:00:00), Dst: Broadcast (ff:ff:ff:ff:ff:ff)
Internet Protocol Version 4, Src: 10.192.19.180, Dst: 127.0.0.1
User Datagram Protocol, Src Port: 0, Dst Port: 54321
INT-Report Protocol Data
  > Telemetry Report Fixed Header
    Encapsulated Header
  > INT Shim Header for TCP/UDP
  > INT metadata Header
  > INT metadata
    > Switch metadata #3
    > Switch metadata #2
      Switch ID: 226
      Ingress port ID: 1
      Egress port ID: 2
      Hop Latency: 25
      0000 0000 ..... = Queue ID: 0
      .... 0000 0000 0000 0000 0000 0000 = Queue Occupancy: 0
      Ingress Timestamp: 1013184455
      Egress Timestamp: 1013184480
      0000 0000 ..... = Queue ID: 0
      .... 0000 0000 0000 0000 0000 0000 = Queue Congestion: 0
    > Switch metadata #1
      Switch ID: 204
      Ingress port ID: 3
      Egress port ID: 1
      Hop Latency: 7
      0000 0000 ..... = Queue ID: 0
      .... 0000 0000 0000 0000 0000 0000 = Queue Occupancy: 0
      Ingress Timestamp: 1014186277
      Egress Timestamp: 1014186284
      0000 0000 ..... = Queue ID: 0
      .... 0000 0000 0000 0000 0000 0000 = Queue Congestion: 0
  > INT Tail Header
  
```

### ③ P4によるINTの実現性の確認 (3/3)

#### ONOSのINT設定画面

Collector IP  Collector Port

**Deploy**

Src Address  Dst Address  Src Port  Dst Port  Protocol

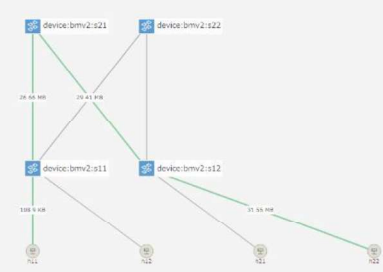
Switch Id
  Port Id
  Hop Latency
  Queue Occupancy
  Ingress Timestamp
  Egress Timestamp
  Egress Port Tx Utilization

**Deploy**

Installed INT Intents (1 total)

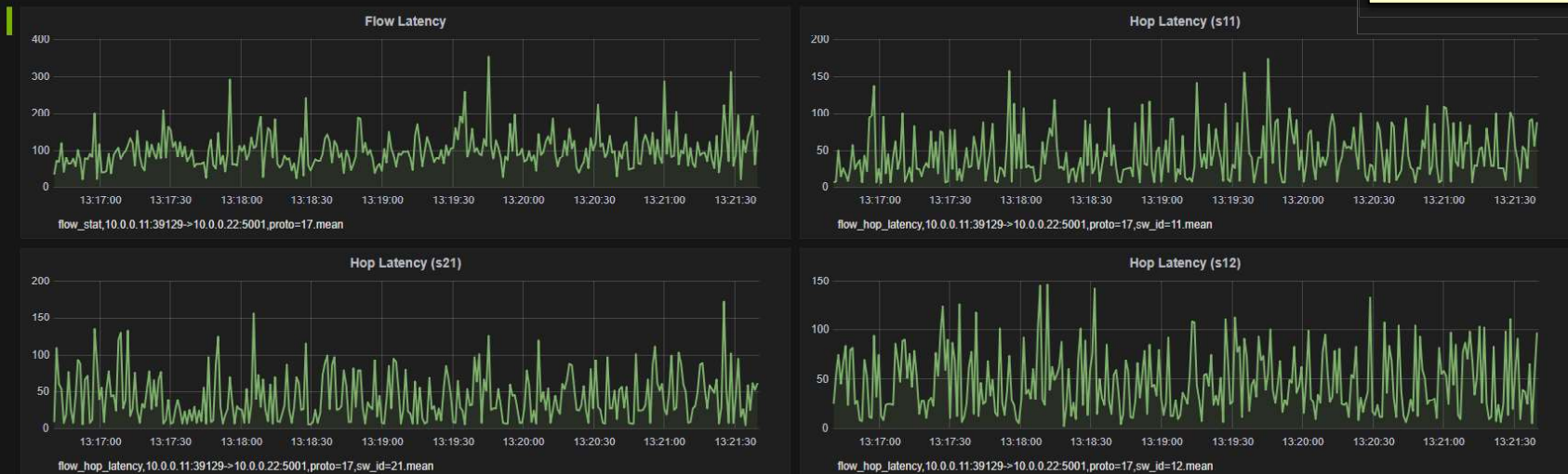
ID	SRC ADDRESS	DST ADDRESS	SRC PORT	DST PORT
1	10.0.0.11/32	10.0.0.22/32	N/A	5001

#### ONOS上のトポロジー図



取得情報の選択が可能

#### Grafana画面



## P4検証を通して、わかったこと

- P4Runtimeの実装が進んできている。
- P4プログラミングは簡単。
  - チュートリアルが整っている。
  - 数百行程度のコーディングで任意の機能がスイッチで実現できる。
- スwitchの環境構築に苦勞
- 現在は、デバイス依存箇所が存在
  - ポート番号、リソース量等のデバイス依存要素を考慮しなければいけない場合あり
  - 検証時、BarefootスイッチはP4\_14にのみ対応 (現在はP4\_16にも対応)

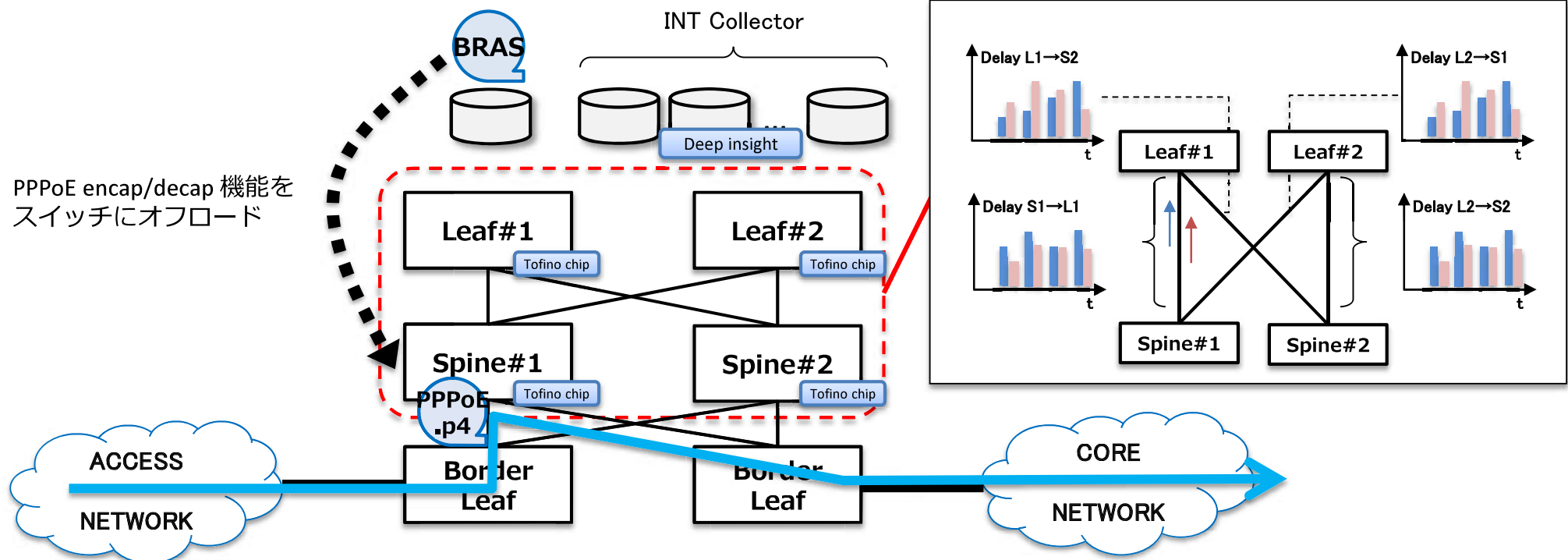
# 今後の取り組み

# PoC @NTT NS Lab.

P4とプログラマブルスイッチを用いたBNGのディスアグリゲーションのPoCをNTT NS研で実施予定(11月)

- ① Disaggregate PPPoE/L2TP Function with P4 program
- ② In-band-network telemetry with P4 program

in-band-network telemetryによって  
各ネットワーク情報を可視化



PPPoE encap/decap 機能を  
スイッチにオフロード

# 今後のスケジュール

- P4/Stratum活用ユースケースの見極めを継続。
  - ✓ 2つのユースケースについて、NTT R&D Forum (2018.11)でPoC予定
  - ✓ 2nd STEPとして、ONOS/Stratum適応について検討。2019.1頃にPoC予定

	FY2018				FY2019
	Apr.	Jun.	Oct.	Jan.	Apr.
<p><b>1<sup>st</sup> STEP</b></p> <p><b>Evaluate P4 Use cases for BNG</b></p> <ul style="list-style-type: none"> <li>➤ Disaggregate PPPoE Function</li> <li>➤ In-band-network telemetry</li> </ul>					
<p><b>2<sup>nd</sup> STEP</b></p> <p><b>Apply P4/stratum for BNG</b></p> <ul style="list-style-type: none"> <li>➤ 1<sup>st</sup> STEP P4 use cases</li> <li>➤ ONOS/Stratum integration</li> </ul>					

▲PoC  
(NTT R&D Forum2018)

▲PoC予定

## まとめ

- P4/Stratumの実装が徐々に進んでいる。
- P4/Stratumを活用のユースケースとして、2つのアプローチを検討。
  - ① VNFのオフロード、②通信品質の可視化
- P4を用いた機能実装は簡単である！