

Disruption Tolerant Networking by Smart Caching

Journal:	<i>International Journal of Communication Systems</i>
Manuscript ID:	IJCS-08-0075.R2
Wiley - Manuscript type:	Research Article
Date Submitted by the Author:	
Complete List of Authors:	Goebbels, Stephan; RWTH Aachen, Faculty 6
Keywords:	Disruption Tolerance, Smart Caching, Intermittent Network Coverage, Extensive Buffering, Pre-fetching



Disruption Tolerant Networking by Smart Caching

Stephan Goebbels
Chair of Communication Networks
Faculty 6, RWTH Aachen University
sgs@comnets.rwth-aachen.de

Abstract: Future mobile radio networks will aim to achieve “broadband access for all”, anywhere. The performance of a radio network vitally depends on the characteristics of the transmission path between the user terminal and the Access Point and the degree of network coverage. In urban areas, full broadband radio coverage is difficult to provide, causing a high variation in the link quality and making broadband services hard to realize. In rural regions, massive deployment costs prevent a full broadband coverage. Most of the time users have to settle for UMTS-like wide area networks. For mobile users accessing services such as video streaming, which require continuous broadband connectivity, it virtually results in intermittent network connectivity. The frequent disruption of the broadband link and its replacement with no or only low-performance connections is a problem that should be addressed. This article introduces a new technique called Smart Caching, which is able to mitigate variations in the network performance so that non-real-time and non-interactive services’ quality is substantially improved.

Smart Caching supports pre-fetching from a server and buffering data at the edge of the core network, in the so called Smart Cache. It transmits data with extremely high speed to be buffered in the mobile terminal when it is in the service range of an Access Point. This allows for the provisioning of data-intensive services even in the case of patchy wireless broadband network coverage and intermittent connectivity.

The performance of the Smart Caching service is evaluated with two different sophisticated queuing models, both based on the Markov Arrival Process. The benefit of the new technique is discussed and dimensioning issues are outlined. Furthermore, a comparison with legacy network setups is given.

Keywords:

Disruption Tolerance, Smart Caching, Intermittent Network Coverage, Extensive Buffering, and Pre-fetching.

1. Introduction

Nowadays ubiquitous telecommunication access is no longer a dream. Mobile users have Internet access wherever they are. But the connectivity is mainly provided by wide area networks like UMTS which offer a peak data rate of only up to 384 kbit/s. In contrast to that in so-called Hot Spot zones broadband Internet is provided by wireless networks like WiMAX. These allow link throughput capacities of several tens of Mbit/s. In addition, the provided data rate always depends on the distance between the access node and the user terminal. Thus, peak data rates are only provided close to the access node – but at the edge of the coverage area the achievable throughput drops to a much lower level. All this results in a highly heterogeneous network coverage. In urban environments, the attenuation of radio waves and shadowing due to obstacles further reduce the coverage range of broadband networks. A mobile user moving in a heterogeneous network such as that shown in the left picture of Figure 1 perceives a continuously changing link performance. Especially if the user tries to access broadband services the data rate in the UMTS network is insufficient. Therefore, the coverage for services with high data rate requirements is virtually reduced to areas of wireless broadband coverage, as shown in the right picture of Figure 1.

In rural areas, the realization of ubiquitous wireless broadband network coverage fails because of high deployment costs. However, for intensively frequented regions such as motorways at least a partial broadband coverage is achievable, e.g., by mounting access nodes at bridges or traffic signs. But all this underscores the fact that in the future, only intermittent broadband connectivity can be provided for mobile users.

1
2
3 For mobile users, frequent disruptions of the broadband network connection have to be
4 compensated to provide sophisticated data-intensive services such as high-quality video
5 streaming. The connection interrupts create a dramatic increase in packet delays compared to
6 legacy networks. Such packet delays have to be countervailed somehow. Nevertheless, for
7 mobile users the delay is inherent to the coverage structure of the network and to user's
8 mobility so that a reduction is impossible. If a packet delay is unavoidable the introduction of
9 delay tolerance to broadband services can solve the problem for adequate service
10 provisioning. A service is only interrupted when the user detects an odd behavior. For
11 example, as long as the playback of a video continues, the service is not regarded as disrupted
12 by the user, although the network connectivity might be lost for a long time. This obviously
13 requires the build-up of a stock of video data in the user terminal. The same holds for push
14 services of personalized (emails, RSS feeds, etc.) or commonly demanded content (newscasts,
15 video streams, etc.). By filling the user terminal with a huge stock of such information, which
16 is of interest to the user, it might be possible that he may not even notice connection
17 interruptions because of consuming the already buffered content. He can, for example, peruse
18 his favorite news page even during connectivity gaps, since the data was pre-fetched in
19 advance. But such an approach requires determining how much data has to be pre-fetched in
20 the end device in order to continue the service even during connectivity gaps. In the case of
21 video streaming, this question is strongly correlated to the packet delay. The longer one has to
22 wait for the next packet, the more data needs to be buffered in the terminal to achieve an
23 uninterrupted service. Therefore, in the following the analysis is focused on the determination
24 of the packet delay under intermittent broadband connectivity. From the packet delay the
25 necessary parameters for the dimensioning of Smart Caching (SC) enabled networks are
26 determined.

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
SC has to take care of an additional challenge: the accumulation of big reservoirs of user data
in the terminal requires a massive transfer of data in periods of good connectivity. Together

1
2
3 with the frequent connectivity losses, it is of vital importance that periods of good link
4 performance have to be optimally used. Then the utilization of the wireless system is brought
5 to its optimum. If it is aimed to fully use the wireless capacity backbone links can have a
6 direct impact on the end-to-end performance.
7
8
9

10
11
12 Therefore SC integrates two fundamental paradigms: first, the optimal utilization of available
13 network resources if broadband connectivity can be provided and second, the build-up of
14 stocks in the end device in the Terminal Buffer (TB) to virtually continue services even in the
15 case of numerous connection interruptions.
16
17
18
19
20

21
22 Section 2 summarizes related work. It also shows that the SC approach is a logical extension
23 of existing research topics. The SC concept is further elaborated in the Section 3. It
24 specifically illustrates how the performance of wireless broadband networks can be
25 optimized. The main focus is on the impact of the backbone on the achievable performance. A
26 performance analysis of SC is presented in Sections 4 – 8. Section 9 concludes.
27
28
29
30
31
32
33

34 35 36 **2. Related Work**

37
38 The use of the Smart Caching approach requires that several aspects have to be taken into
39 account. The cooperation of SC with the existing network infrastructure has to be considered.
40
41 Cache replacement strategies as well as content and user path, respectively next visited cell,
42 prediction are strongly related to the SC technique.
43
44
45
46
47

48 A problem associated with temporarily storing user data is that it might interfere with network
49 protocols like TCP as TCP is an end-to-end connection oriented protocol. Therefore it is
50 necessary to decouple the connection into two subsections – one between the server and
51 Smart Cache (SCache) and another between the SCache and Terminal Buffer. For TCP a
52 solution, called SNOOP, was already presented in [1]. It was developed for the use of TCP
53 over wireless links, which suffers from similar problems. Nevertheless the SNOOP approach
54 fails to further extend the separation between wired and wireless links. The introduction of a
55
56
57
58
59
60

1
2
3 buffer like proposed in SC allows a much better performance enhancement then the simple
4 interception of the end-to-end protocol connection known from SNOOP.
5
6

7
8 In [2] a special extended RTP protocol is proposed, which is focused on streaming traffic
9 including caching devices and allows reliable data delivery through caches. It does not require
10 any changes in the end user software so that it is transparent to it. Video streaming,
11 investigated in the later analysis, can be supported by such protocols.
12
13
14

15
16
17 Although due to the clustering of APs it is not required to predict the next visited cell in most
18 of the cases it can further boost the performance of SC. In [3] a good survey about next cell
19 prediction techniques for wireless networks is given. It uses long term log files of pedestrian
20 user mobility as training sequences for different analytical models. It is shown that even with
21 low complexity models a good mobility prediction is possible. Hence, further enhancement of
22 the SC approach is possible. However, the long initial monitoring intervals might cause a
23 lower performance at the deployment phase.
24
25
26
27
28
29
30
31
32

33
34 The pre-fetching of data in the end device makes it necessary that the content of the cache is
35 continuously replaced by more recent information. In [4] updating strategies for caches are
36 discussed. Different strategies are compared which take into account content recency or
37 access frequency. In [5] personal profiles of web surfing behavior are used to predict the next
38 accessed content. Both approaches allow a most accurate pre-fetching of data which is of
39 further benefit for SC.
40
41
42
43
44
45
46
47

48 The SC approach is highly related to Delay Tolerant Networking (DTN). SC also tries to
49 handle communication in network setups with long packet delays. DTN, started from
50 interplanetary communication, noted that terrestrial communication also suffers under specific
51 circumstances from similar problems. DTN tries to overcome connectivity interruptions by
52 using the terminals of other mobile users for data transport between the access node and the
53 destination terminal. SC instead follows the direct way. The mobile user himself reaches the
54 next coverage zone and accesses the data there. The intermittent connectivity is compensated
55
56
57
58
59
60

1
2
3 by extensive buffering. In both approaches, long packet delays are expected. But the success
4
5 of DTN already shows that in the future such long delays are acceptable for mobile users.

6
7 While sensor networks can accept longer packet delays without for human users
8
9 compensation techniques have to be found – extensive buffering in the case of SC.

10
11 In [6] the general architecture, the design, and the state of the art of DTNs are outlined. The
12
13 article does not consider the performance but focuses on the protocol and design aspects. The
14
15 proposed protocols can be used to make SC working. In [7] the store and forward principle of
16
17 DTNs is discussed. The analysis in this article focuses on regions with rudimentary network
18
19 coverage which should be enhanced by an intermittent broadband connectivity. Such areas
20
21 should be supplied by so called data mules which physically transport the data. For example,
22
23 consider a small village connected via one wireless link for basic access provisioning, and
24
25 additionally a frequently operating bus piggybacks the data between an access gateway in the
26
27 next bigger city and the village. From this approach the organization of the store and forward
28
29 principle can be applied for the current concept.
30
31
32
33
34

35
36 In [8] and [9] performance results of DTN are presented. Both articles show that packet
37
38 delays of several tens of seconds are bearable. In the latter article it is shown that the buffer
39
40 size has to be taken into account as it is a finite resource and can impact the performance of
41
42 the approach. This shows that the later analysis of the TB size considers most relevant
43
44 parameters. The impact of the mobility model on the performance evaluation of intermittent
45
46 networks is self-evident. As used in the later analysis in [10] a Manhattan like scenario is
47
48 taken to evaluate the performance of SC. Although DTN and SC are much related the former
49
50 simply tries to exploit the delay resistance of specific services, e.g., data gathering within
51
52 sensor networks. SC goes further and tries to mobilize (broadband) services which are not
53
54 supportable in legacy networks due to their delay bounds. Therefore not the service is adapted
55
56 to the network but the networks adapts to the needs of the service.
57
58
59
60

1
2
3 The idea to transfer data only in regions where high bit rates can be provided was already
4 followed by [11]. This approach also inspired the later presented motorway scenario. In the
5 article only general ideas and concepts are outlined so that in the following analysis the focus
6 is put on the achievable improvement in system performance and network capacity.
7
8
9
10
11

12 13 14 15 **3. The Smart Caching Concept**

16
17 The architecture of Smart Caching entails two new network devices: the Smart Cache
18 (SCache) and the Terminal Buffer (TB), as sketched in Figure 2. Caching in the SCache, as
19 described in this article, serves for the reduction or even removal of the “inner resistance”
20 (resulting from flow control algorithms or from capacity constraints) of a fixed
21 telecommunication network, as seen by a user terminal. To achieve this, data from a server is
22 pre-fetched and stored in the SCache so that the mobile terminal can access user data with the
23 maximum transmission speed its wireless link is able to support. So instead of fine tuning
24 transmission protocols to slightly increase the throughput SC introduces a completely new
25 network setup so that the end-to-end data rate can be brought to its optimum.
26
27
28
29
30
31
32
33
34
35
36
37

38 Which impact SC can have on the data rate is shown in Figure 3. The graphs do not show real
39 world measurements but instead the general characteristics which wireless and wired
40 networks reveal for mobile users. In each graph the data rate of the wireless network, the
41 backbone network and the resulting end-to-end data rate is displayed for a terminal moving
42 from left to right through the center of the cell. The data rate of the wireless network as seen
43 by the user terminal is of pyramid shape. Periods of broadband wireless network coverage
44 may be followed by periods without any broadband connectivity (no connectivity or only
45 service by a cellular wide-area network). Accordingly the radio data rate is set (close) to zero
46 between two phases of AP connectivity in the example shown. The percentage of time where
47 network connectivity is provided is denoted by the coverage ratio (cov).
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3
4
5
6
7
8
9
10
11
12
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

Contrary to wireless networks in the backbone the achievable data rate is relatively constant. This data rate is in the following called the Backbone Limit (BL) as it might prevent the wireless network from achieving its best performance. The backbone data rate will usually not achieve the peak rates of the wireless network. In legacy networks, the resulting end-to-end data rate is always the minimum rate of all involved links, here the wireless and the backbone connection. This behavior “cuts of” a substantial part of the possible performance of the wireless network, see upper graph of Figure 3.

By buffering data in the AP, the initial phase of low radio data rates might be used to build up a stock, which then is forwarded when enough capacity in the wireless network becomes available (second graph of Figure 3). But this has only a minor impact on the overall performance.

Only by caching packets which arrive during the gap period it is possible to gather a sufficient amount of data which makes it possible to substantially utilize the otherwise unused capacity of the wireless network. To make use of the cached data, it is necessary that it is stored close to the next visited AP. Since the prediction of which AP will be visited next is usually difficult, all APs of a closer proximity have to be clustered and connected to one SCache.

Such an operation involving massive buffering is unsuitable for interactive communication services like speech and video conference owing to strict delay bounds specified for these services' data packets. But these applications may benefit from the service provided by a wide area radio network, since they are not bound to wireless broadband networks.

Non-real-time services like unidirectional video streaming, FTP, email push/pull, and similar services tolerate high packet delays. SC exploits this property by pre-fetching contents and buffering it to transfer it at maximum possible speed to the user terminal. The corresponding storage is the second new network device introduced by SC the TB.

From the perspective of communication partners, the application server and the service user, the end-to-end connection is enhanced with two buffers. One is the SCache in front of the

1
2
3 wireless link and one is the TB directly behind the wireless link. These two buffers together
4
5 create a black box which is fed with data from the server, see Figure 4. This data is consumed
6
7 in the end device by the user.
8
9

10 The TB is a storage which keeps a substantial data reservoir in the end device. Owing to the
11
12 connection interruptions, a mobile user perceives the delivery of data to be much delayed. In
13
14 Figure 5 it is shown how a coverage gap interrupts the data delivery and the packets are
15
16 delayed. The time axis is oriented in the negative y direction. Packets are reflected by
17
18 horizontal arrows. SCache and AP are regarded as colocated since the data exchange between
19
20 these nodes is so fast that it has no impact on the performance of the end-to-end connection.
21
22 Packets arrived in the SCache/AP are usually forwarded to the terminal with a short delay. In
23
24 the terminal the data packets are buffered in the TB until the included content is displayed on
25
26 the screen.
27
28
29

30
31 But in case of a coverage gap, the packets have to be buffered until the connection is regained.
32
33 This is done in the SCache. For video streaming, therefore, it is necessary to determine how
34
35 big the reservoir in the TB has to be to compensate the packet delays caused by the
36
37 connection interruptions.
38
39

40
41 The right hand side of Figure 5 illustrates that the delay can be compensated by shifting the
42
43 playback of the video in time. This time shift has to be big enough to compensate the packet
44
45 delay. In the case of a coverage gap, the data in the TB is consumed but the new packets
46
47 should arrive early enough to continue the playback before the TB is drained. At the outset it
48
49 is only necessary to store an amount of data in the TB that corresponds to the required time
50
51 shift before the actual playback starts. Therefore SC is the first technique which allows
52
53 mobilization of broadband services which are currently limited to, so called, Hot Spot zones.
54
55 The network configuration is changed in such a way that service interrupts are so much
56
57 reduced that the user perceives a virtually continuous broadband connection. Since the time
58
59
60

1
2
3 shift vitally depends on the packet delay, in the following analyses a special focus is put on
4
5 the determination of the end-to-end packet delay.
6
7

8 The SCache works as a mediator between the wired and the wireless world by introducing a
9
10 buffer between them so that the end-to-end data flow is separated into two subsections. The
11
12 SCache terminates the first subsection commencing at the service provider. Concurrently, it is
13
14 the starting point of the second subsection which terminates in the TB.
15
16

17 The actual communication partners are the service provider offering, e.g., a video stream, and
18
19 the user terminal by which the service is requested. The SCache resides in the middle between
20
21 the two communication partners. For the service provider the SCache adopts the role of the
22
23 data flow end point. At this point the end-to-end connection between the service provider and
24
25 the client is interrupted. The ingress data leaves the transport layer and moves up to a data
26
27 buffer where it is cached until it can be finally forwarded to the user terminal. The first
28
29 subsection exclusively consists out of wired Internet connections. Therefore TCP or similar
30
31 protocols are suitable for the transfer of data in this subsection.
32
33
34

35 Pre-fetching and buffering is only reasonable if the data actually reaches its final destination,
36
37 the user terminal. Data which is pre-fetched in the SCache but does not arrive at the end user
38
39 wastes network resources. To locate the SCache directly in an AP, like in the SNOOP attempt
40
41 [1], would imply that all pre-fetched data is useless, if the user leaves the coverage range of
42
43 this AP and never returns. Making use of already buffered data requires that it must be
44
45 available also from other access nodes which the user possibly roams to in the near future.
46
47
48

49 If all APs of a certain area are clustered and connected to one SCache it is possible to benefit
50
51 from the cached data as long as the user terminal is supplied by one of these APs. Starting a
52
53 communication session in cell A the data stream is routed through the SCache which is
54
55 responsible for the cluster that AP cell is assigned to. If the data rate on the wireless link
56
57 decreases because either the user moves away from the AP or completely leaves the cell, the
58
59 data transmission on the first subsection continues. But as no forwarding through the second
60

1
2
3 subsection is possible, the data is buffered in the SCache. Entering a new AP cell **B** the
4
5 communication session will be re-established. If the new AP is included in the same cluster as
6
7 cell **A** it is sufficient to reopen the connection between SCache and user terminal while the
8
9 first subsection between service provider and SCache remains unchanged.
10
11

12 The analysis in the Sections 4 – 8 proves that SC-supported networks allow continuous video
13
14 streaming service even under intermittent network coverage with long connectivity gaps.
15
16

17 Furthermore, the advantage of the new approach over existing techniques is shown.
18
19

20 21 22 **4. Urban Application Scenario** 23

24 To analyze the performance of Smart Caching, especially its ability to compensate
25
26 connectivity gaps in urban environments, the well-known UMTS 30.03 Manhattan scenario is
27
28 chosen. All the APs in the scenario are connected to and served by one Smart Cache. The
29
30 street grid consists of 200m-sized blocks and 30m-wide streets. WiMAX APs are placed on
31
32 every second crossover. They are operated with a transmission power of 100 mW. WiMAX
33
34 uses different modulation and coding schemes (PHY modes) to adapt the radio transmission
35
36 to the current radio link quality. Each of them is assigned with a minimum signal quality
37
38 necessary to decode the information and an available data rate (network parameters for
39
40 WiMAX are taken from [12]). The pathloss between sender and receiver is given by an
41
42 adapted free space model with an attenuation factor γ of 3 (2 stands for free space – up to 5
43
44 for indoor). Transmission power, pathloss and minimum signal level allow the mapping of a
45
46 PHY mode to each terminal position in the scenario. As shown in Figure 6 the streets are
47
48 divided in annulus shaped zones served by a specific PHY mode and zones without WiMAX
49
50 connectivity. The network coverage ratio (for streets) is given by 82%.
51
52
53
54
55
56

57 To get representative input parameters for the later analysis in a mobility simulation the
58
59 behavior of users in such a scenario is derived. Pedestrian users are traversing the streets of
60
the scenario with a velocity equally distributed between 0.5 and 1.5 m/s. Their mobility is

1
2
3 subject to an adapted Brownian motion which includes velocity updates and direction changes
4
5 of up to 45 degrees on average every 10 seconds. The outcome (the user's pathway) of the
6
7 first 10000 seconds of such a simulation is shown in Figure 6 (right). Out of much longer
8
9 simulation runs reliable values for the average residence time in each PHY mode zone and the
10
11 transition probability of users between different PHY modes are gained. Additionally, a large
12
13 sample set of gap period durations can be derived.
14
15

16
17 For later use, a pure sample set is insufficient; thus the complete Cumulated Distribution
18
19 Function (CDF) of the gap duration sample set is approximated with the support of the EM
20
21 algorithm [13]. The outcome is a phase-type distribution which provides a closed-form of the
22
23 CDF. The results for the gap duration distribution function and its approximation are shown in
24
25 Figure 7.
26
27

28
29 Furthermore, the parameters of the different PHY mode zones can be employed to set up
30
31 models for arrival and service processes, see Section 5. Out of the coverage ratios of each
32
33 PHY mode and the corresponding data rates, it is possible to calculate the average data rate a
34
35 user perceives. While for the backbone limited setups the performance of WiMAX is cut to a
36
37 certain boundary the SC-enabled setups allow an optimal utilization of the network resources.
38
39 In Figure 8 the available average data rate a user perceives for a test service depending on the
40
41 Backbone Limit with the capacity share as a parameter is shown. The dotted red lines display
42
43 SC-enabled setups and the solid black lines display the legacy setups without SC. The
44
45 capacity share reflects the portion of radio resources that the test service can use and is not
46
47 occupied by the other services (irrespective of whether they use SC or not). If network
48
49 resources are bounded to other services the wireless network can only be used partially for the
50
51 data transport of the test service. This can be reflected by a reduction of the available data
52
53 rate. For example, if 75% of the resources are occupied by other services (corresponds to 25%
54
55 capacity share), the data rate available for the SC service is automatically cut to one-fourth.
56
57
58
59
60

1
2
3 Assuming that the mobile user streams and watches a high quality video an average data rate
4 of 3.37 Mbit/s can be assumed (MPEG4 Codec, see [14] page 108). Considering the curves of
5
6
7
8 Figure 8, it can be seen that the advantage of the SC-enabled network for such a data rate is
9
10 not directly evident. Especially for large capacity shares, the difference between the two
11
12 curves (red dashes and solid black) is marginal for a wide range of the BL. But for a capacity
13
14 share of only 20%, the advantage increases. While without SC the BL must be above 7.6
15
16 Mbit/s this value can be decreased to the actual streaming rate of 3.37 Mbit/s which is less
17
18 than one half. Later on it will be shown that the advantage of SC is that the achievable data
19
20 rate drastically increases compared to the legacy case if the coverage ratio gets lower
21
22 (between 20 and 50%).
23
24
25

26
27 A capacity share of 20% implies that with full capacity at least five mobile video users can be
28
29 served simultaneously per cell. For approximately 0.022 km² covered street area per WiMAX
30
31 cell it means one user per 4000 m² which in a densely populated urban environment is not an
32
33 unlikely value. Due to statistical multiplexing of users residing in a coverage zone or being
34
35 outside of it the number of acceptable users could be further increased. Thus, depending on
36
37 the BL the service might not be realizable without SC as the average data rate would be
38
39 insufficient. But however there still exists the problem of the duration the packets spend in the
40
41 SCache. Therefore in the following a sophisticated queuing model is used to develop the
42
43 delay.
44
45
46
47
48
49

50 **5. MMAP/G/I Queuing System**

51
52 As stated earlier, the packet delay is essential for the dimensioning and performance
53
54 evaluation of Smart Caching. Modern queuing theory allows a detailed modeling of different
55
56 arrival streams competing for the limited radio resources. As a direct output the packet delay
57
58 of the different streams becomes available.
59
60

1
2
3 The queue itself can be seen as a combination of the Smart Cache and Access Point, compare
4 Figure 2. The SCache buffers the packets until they can be forwarded via the wireless link to
5 the user terminal. The wireless link is the server of the model which works off the queued
6 packets. The Terminal Buffer and the actual display of the user terminal can be regarded as a
7 second queuing system. The arrival process represents packets which are transferred via the
8 wireless link. The TB is reflected by the queue. Packets are worked off by the display when
9 the corresponding packet content is shown on the screen. The next analysis is focused on the
10 first queuing system. The second system is discussed later on and the results of both analyses
11 are compared.
12
13
14
15
16
17
18
19
20
21
22
23

24 To show the benefit the employment of SC implicates the approach has to be compared with
25 legacy network setups. Without SC the Backbone Limit reduces the achievable data rates
26 between server and client. By adapting the service process of the queuing system both
27 approaches can be compared.
28
29
30
31
32
33

34 In [15] new results of queuing theory are introduced, which allow the modeling of queuing
35 systems with different incoming packet streams, each of them represented by a Markov
36 Arrival Process (MAP) and can be subject to a separate service time distribution. The queuing
37 delay of each arrival stream can be calculated separately. Since it can be distinguished
38 between the arrival streams the process is called marked MAP (MMAP) and the resulting
39 queuing system has the notation $MMAP(i)/G(i)/I$. This allows a very detailed modeling and
40 analysis of the investigated scenario setup.
41
42
43
44
45
46
47
48
49

50 A MAP can be described by a Markov chain which allows transitions between all states. But
51 contrary to a normal Markov chain two types of state transitions can occur. The first type is
52 the transition between two internal states which is not noticeable from the outside and the
53 other is a state transition with a concurrent arrival. The transition rates of the two types are
54 summarized in the matrices D_0 and D_1 . D_0 contains simple state transitions; D_1 covers
55 transitions with arrivals. For each arrival process i a separate pair of matrices $D_{0,i}$ and $D_{1,i}$
56
57
58
59
60

exists. The integration of the different streams can be done by composing the Kronecker sum of the different matrices $D_{0,i}$ and $D_{1,i}$, which is described in [16]. At the end there is only one matrix D_0 and i different matrices D'_i according to the separate traffic streams. The dimension of the resulting matrices arises from the product of the size of the input matrices $D_{0,i}$.

Two types of arrival processes have to be considered. First, a pure Poisson process ($D_{0,i}=-\lambda$ and $D_{1,i}=\lambda$, λ = average arrival rate) to reflect background traffic which consumes network resources that are unavailable for the SC-supported traffic. This can be used to vary the capacity share. And secondly, the actual video streaming process which is supported by SC has to be modeled.

Already in [17] it was proposed to model the Variable Bit Rate (VBR) video traffic by using a certain number M of so called mini sources. To do so the video streaming rate is quantized in chunks of size λ_q . Each quantization step is modeled by a mini source which produces data with an average rate of λ_q . Depending on how many mini sources are active, the accumulated data rate varies as that by VBR traffic. The number of active sources is steered by an underlying birth death process as shown in Figure 9.

The above modeling of the video streaming process corresponds to a Markov Modulated Poisson Process (MMPP). An underlying Markov chain steers the state changes and to each state a data rate is assigned. A MMPP consists of two matrices Γ and Λ . The first contains the state transitions and the second the arrival rates per state. The translation to a Markov Arrival Process is achieved by

$$\begin{aligned} D_0 &= \Gamma - \Lambda, \\ D_1 &= \Lambda. \end{aligned} \tag{1}$$

The transition rates α and β as well as λ_q of Figure 9 can be matched to the parameters of the video streaming process. In [17] the average streaming rate $E[\lambda]$, the variance $C(0)$ and the auto-covariance function $C(a,\tau)$ are used for that. With the following equations the necessary parameters can be deduced.

$$\alpha = a \left(1 + \frac{E^2[\lambda]}{MC(0)} \right)$$

$$\beta = a - \alpha \quad (2)$$

$$\lambda_q = \frac{C(0)}{E[\lambda]} + \frac{E[\lambda]}{M}$$

The actual parameters for MPEG4-coded video streaming are taken from [14], see Table 1.

The service process of the queuing model has to reflect the different PHY modes by which a mobile user is served. To reflect the different data rates of each PHY mode in the queuing model a hyperexponential service time distribution is taken, see Figure 10. According to path probabilities p_i packets are processed with different service rates μ_i . Both parameters can be derived out of the probabilities that a user resides in PHY mode area i and the corresponding data rate r_i of the WiMAX system. These values are provided by the mobility simulation and the wireless network characteristics.

To fully understand the benefit of SC, it is important to compare its performance with the legacy case. To emulate the “cutoff” behavior of legacy network setups as shown in the upper graph of Figure 3 the service rates μ_i have to be reduced accordingly.

Up to now the service process does not cover areas without connectivity. As the transmission rate is zero in such regions, it is not possible to directly include it in the above model. If a customer traverses a coverage gap the incoming packets of the, e.g., video stream are buffered in the SCache. Due to this buffering the end-to-end delay of packets is drastically increased.

Which influence the gap has on the actual packet delay is depicted in Figure 11. The scene is separated in three periods. In the first period, denoted by x , the customer leaves a coverage zone and the SCache starts the buffering so that the fill level increases. After re-entering a new coverage zone the incoming traffic is buffered as there are still earlier arrived packets left in the cache which have to be served first. Since the arrival rate for new packets is smaller than the transmission rate of the air interface the buffer starts to get drained and the fill level decreases. The period until the whole buffer is emptied is called y . And finally the period z

1
2
3 just denotes the normal operation of the wireless network. Packets which arrive have to wait a
4
5 short period until radio resources become available to be transmitted.
6
7

8 In the lower part of Figure 11 the additional delay depending on the arrival time of the packet
9
10 in the SCache is shown. It is caused by the gap in the coverage. The maximum delay is
11
12 suffered from the first packet which arrives after the user terminal has left the coverage zone.
13
14 It is delayed for the whole period of no coverage and is instantly transmitted when the next
15
16 coverage zone is reached. Therefore the additional delay vitally depends on the duration of the
17
18 gap period. The packet delay linearly decreases until the point is reached where the complete
19
20 SCache is drained and the operation migrates to the normal behavior. For later evaluations it
21
22 is necessary to determine the ratio of packets which are affected by the gap and the portion
23
24 which is transferred without additional delay. Clearly the first portion is reflected by
25
26 $(x+y)/(x+y+z)$.
27
28
29
30

31 The amount of traffic which has been aggregated in the SCache during a gap is given by
32
33 $x\lambda_{video}$, where λ_{video} is the arrival rate of the video process. During the y period this amount of
34
35 stored data together with the currently arriving traffic has to be worked off. The current traffic
36
37 loads the system with the utilization ρ . The utilization is defined as the quotient out of the
38
39 average arrival rate λ of a queuing system and the average service rate μ . The arrival rate is
40
41 predefined by the data rate of the regarded user services (e.g., video streaming rate). The
42
43 service rate in the regarded system model reflects the ability of the end-to-end connection to
44
45 transport the data packets. This is influenced by two aspects the capacity of the wireless link
46
47 but as well the BL, which might prevent an optimal performance of the radio connection. The
48
49 capacity share which is available to work off the aggregated data heap is $(1-\rho)$. All this results
50
51 to
52
53
54
55
56

$$x\lambda_{video} = y(1-\rho)\mu_{video} \quad (3)$$

57 where μ_{video} is the video service rate. Together with the fact that the coverage ratio is given by
58
59 the quotient out of x and the overall size $x+y+z$, it can be concluded that
60

$$\frac{x+y}{x+y+z} = (1-\text{cov})\left(1 + \frac{\rho_{\text{video}}}{1-\rho}\right). \quad (4)$$

The additional delay the packets perceive due to the gaps in the connectivity is equally distributed between gap duration and zero. Since this delay is independent from the waiting time in the queuing model both values can be summed up. For their probability densities this implies a convolution. Since the outcome of the *MMAP/G/I* queuing system is the Laplace Stieltjes Transform (LST) of the waiting time distribution it is natural to perform the convolution by a simple multiplication of the LSTs. Therefore it is also necessary to get a closed form expression of the LST of the gap duration. As packets during the period z out of Figure 11 are not affected by the gap it is necessary to separate the distribution of the additional delay in two parts. With the probability $z/(x+y+z)$ the additional delay is zero and with the probability $(x+y)/(x+y+z)$ an additional delay exists. The first part can be reflected by a Dirac impulse in the origin $\delta(t)$. For the second part the probability density function (pdf) of the additional delay $pdf(t)$ (compare the proof at the end of the article) is required. $pdf(t)$ and the resulting $pdf(\text{delay}=t)$ for all packets are given in Equation (5).

$$pdf(t) = \int_t^{\infty} \frac{p(x)}{E[x]} dx = \frac{1}{E[t]} (1 - cdf(t)) \quad (5)$$

$$pdf(\text{delay}=t) = \frac{z}{x+y+z} \delta(t) + \frac{x+y}{x+y+z} pdf(t)$$

What is still missing for a final solution of the packet waiting time is the waiting time expression of the *MMAP/G/I* system. The LST of the waiting time distribution per packet type i is given by

$$W_k(s) = \frac{V(s)D_k}{\lambda_k}$$

$$V(s) = (1-\rho) \mathbf{sg}(sI + D0 + D(s))^{-1} \quad (6)$$

$$D(s) = \sum_k D_k(s) = \sum_k D_k H_k(s) = \sum_k \int_0^{\infty} e^{-st} D_k(t) dt$$

where λ_i is the arrival rate of type i packets and ρ is the above mentioned utilization of the system. The vector \mathbf{g} is the stationary vector of the matrix Q ($\mathbf{g}Q=I$ and $\mathbf{g}I=I$) which is given by the recursive formula

$$Q = D_0 + \int_0^{\infty} D(x)e^{Qx} dx. \quad (7)$$

$D(x)$ is similarly defined to $D(s)$ but instead of using the LST the probability densities have to be used. The matrix Q can be iteratively calculated by starting with $Q=D_0$ and continuously substituting it in Equation (7) until the differences between concurrent results gets small enough.

With all this preparation it is possible to derive the CDF of the video streaming packet delay as depicted in Figure 12 (please note that the y-axis ranges from 0.5 to 1.0) for the transport between server and TB. The 95th percentile of the packet delay is a good measure for the required time shift. The scenario is still the urban environment as introduced above. The capacity share is reduced by varying the arrival rate of the additional Poisson process. From the curves as well as from the 95th percentile of the waiting time (listed in the legend) it can be seen that the capacity share has an impact on the waiting time. However, even the provisioning of the full network capacity for one video streaming process can not reduce the waiting time below a certain limit.

If all wireless network resources are reserved for the observed video streaming service (100% capacity share) the 95th percentile of the packet delay and therefore the required time shift accounts to 102 seconds. This means that 102 seconds of video data have to be buffered in the TB before the video playback should start. Under this condition the risk of a playback interruption is much reduced. Later analyses will show that it is around 5%. Since at the beginning of a streaming process the transfer rate can be drastically increased compared to the actual streaming rate it is possible to download the time shift reservoir faster than the 102 seconds. If the wireless network provides in the beginning a data rate which is twice the streaming rate it implies that the required data can be downloaded in around 50 seconds. Although this value is higher than used from legacy streaming services SC allows the service provisioning of high quality video streaming under intermittent network connectivity, which is definitely not possible with legacy setups.

1
2
3 Even for a capacity share of 25% the SC-enabled scenario setup still performs satisfactorily.
4
5 The 95th percentile of the packet delay only increases to 142 seconds so that still a reasonable
6
7 time shift is reached. In the throughput analysis above it was stated that a legacy network
8
9 setup under this condition might get into problems. Therefore the performance of the SC
10
11 approach is compared with a legacy setup. For the legacy case the service rate of the WiMAX
12
13 APs is bounded by the BL. For different BLs the CDF is shown in Figure 13 (capacity share =
14
15 25%). The curves get even lower and the 95th percentile further increases.
16
17
18

19
20 Although the BL has some impact on the packet waiting time the changes are not dramatic.
21
22 The resulting time shift is not influenced much so that the service can still be supported. Only
23
24 the initial download phase is increased by 80% compared to a full capacity share and no BL.
25
26 But the ability of SC to compensate the BL has another much more important relevance. It is
27
28 an enabler for broadband services. Without SC, many services cannot be supported at all.
29
30

31
32 If the BL becomes less than 8 Mbit/s in the given scenario the packet delay goes to infinity as
33
34 the service can no longer be supported by the network. The reason is that the available
35
36 average data rate drops below the playback rate of the video.
37
38

39
40 Nevertheless, even for the SC-enabled network setup the problem on how the delay can be
41
42 handled exists. For applications like video streaming the delay can be compensated by
43
44 extensively buffering video data in the end device. If data storage is provided in the end
45
46 device (namely the Terminal Buffer) the playback of the video can continue when the mobile
47
48 user leaves the zone of wireless network connectivity – as long as there is still buffered video
49
50 data at hand. The question is, how big the TB has to be? And secondly, what is the required
51
52 initial fill level before the playback of the video can start? A reference value for the initial fill
53
54 level of the buffer could be again the 95 percentile of the packet waiting time. It has to be
55
56 guaranteed that new packets arrive early enough to continue the video playback before the
57
58 SCache is drained. If an amount of video data is stored in the end device which corresponds to
59
60

1
2
3 the 95th percentile of the CDF it means that with 95% probability the next packet will arrive
4
5 within that time limit.
6
7

10 **6. MAP/G/I/N Queuing System**

11
12 Now, the second queuing system is regarded. The input process reflects the amount of data
13 transferred via the wireless link. As with Smart Caching always enough data is provided in
14 the Smart Cache the capacity of the wireless link is always fully utilized. Therefore the arrival
15 rate basically reflects the available wireless data rate. The Terminal Buffer is reflected by the
16 queue and the service process models the playback of the video. To allow good service, the
17 TB size has to be determined. And to identify the necessary TB size, the SC-enabled system is
18 considered from the end device's perspective. The TB tries to build up a stock of video data
19 that can be consumed during phases of non-existent or insufficient network connectivity to
20 continue the playback of the video stream. To achieve this goal, it is necessary to fill up the
21 buffer in periods of perfect network connectivity to its maximum. The question is, which
22 buffer size is required to provide a predefined Quality of Service (QoS) level? An obvious
23 QoS criterion is the probability that a video stops, or in other words, the TB is drained.
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40

41 Such an end device together with a limited TB can be modeled by a queuing system with
42 finite queue size. The *MAP/G/I/N* system presented in [18] allows the thorough modeling of
43 the arrival process as well as the monitoring of the queue size and the idle probability of the
44 queue.
45
46
47
48
49

50 The key factor for the analysis of such a queuing system is given by the transition probability
51 matrix P of the counting process which monitors the number of packets in the finite queue.
52
53

54 The matrix is given by
55
56
57
58
59
60

$$P = \begin{pmatrix} \bar{D}A_0 & \bar{D}A_1 & \cdots & \bar{D}A_{N-1} & \bar{D}A'_N \\ A_0 & A_1 & \cdots & A_{N-1} & A'_N \\ 0 & A_0 & \cdots & A_{N-2} & A'_{N-1} \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ 0 & 0 & \cdots & A_0 & A'_1 \end{pmatrix}. \quad (8)$$

For the matrices A_i and A'_i in [18] easy calculation formulas are provide in case the service time distribution is subject to a phase type distribution. Since this type includes a wide area of possible distributions the modeling flexibility is not significantly restricted by that condition.

The matrix \dot{D} can be derived from the arrival process. \dot{D} results to $(-D_0)^{-1} D_1$.

The steady state vector \mathbf{x}' of the matrix P ($\mathbf{x}'P = \mathbf{x}'$ and $\mathbf{x}'\mathbf{1} = 1$) provides the queue size distribution at departure epochs. Departure epochs are events when a packet leaves the queue and is handed over to the server. The whole vector can be separated into subvectors. The entries of one of these subvectors represent situations with the same buffer fill size but different states of the arrival process; $\mathbf{x}' = (\mathbf{x}'_0, \dots, \mathbf{x}'_N)$.

The queue size distribution at arbitrary times ($\mathbf{x}_0, \dots, \mathbf{x}_N$) and the probability that the whole system is empty (including the server) $\tilde{\mathbf{x}}$ can be derived from the following formulas.

$$\begin{aligned} \tilde{\mathbf{x}} &= \frac{1}{E^*} \mathbf{x}'_0 \bar{D} (-D_0)^{-1} \\ \mathbf{x}_0 &= \left(\frac{1}{E^*} (\mathbf{x}'_0 - \mathbf{x}'_1) - \tilde{\mathbf{x}} D_1 \right) D_0^{-1} \\ \mathbf{x}_n &= \left(\frac{1}{E^*} (\mathbf{x}'_n - \mathbf{x}'_{n+1}) - \mathbf{x}_{n-1} D_1 \right) D_0^{-1}, \quad 1 \leq n \leq N-1 \\ \mathbf{x}_N &= \left(\frac{1}{E^*} \mathbf{x}'_N - \mathbf{x}_{N-1} D_1 \right) D_0^{-1}, \text{ with} \\ E^* &= \frac{1}{\mu} + \mathbf{x}'_0 (-D_0)^{-1} \mathbf{1} \end{aligned} \quad (9)$$

The most interesting vector is $\tilde{\mathbf{x}}$ as it contains the probability the system is idle which means the playback of a video has stopped ($P_{idle} = \tilde{\mathbf{x}}\mathbf{1}$).

The size N of the finite buffer has a direct impact on the size of the matrix P . Due to memory and processing speed restrictions the number of entries in the matrix should be limited. A buffer size of $N > 500$ storage units is not recommended. If one data packet would be used as a

1
2
3 storage unit the queue could only hold up to 500 packets. This is much too less for the
4
5 considered approach. Therefore it is necessary to quantize the amount of buffered data. The
6
7 most significant expression for the buffer size is the number of seconds the video can be
8
9 continued with the data stored in the buffer. Therefore the quantization is done on data chunks
10
11 which are enough to playback one second of the video stream (in the case of 3.37 Mbit/s
12
13 playback rate one second corresponds to 0.42 MByte of video data).
14
15

16
17 For the service time of the queuing system a negative exponential distribution is chosen. Since
18
19 it only reflects the consumption of the buffer by the video playback, no more specific
20
21 distribution is taken here. Owing to the quantization of the buffer content, the rate μ of the
22
23 service process is simply given by $1/s$. This means that per second, exactly one chunk of
24
25 buffered data is consumed by the playback of the video.
26
27

28
29 Contrary to the service process, the arrival process of the second queuing system requires
30
31 much more modeling effort. The arrival process is not so much influenced by the actual traffic
32
33 stream it carries, but to a large extent by the data rate that the wireless network can provide.
34
35 Due to the SC approach, always the entire available network resources are consumed. For a
36
37 mobile user traversing an urban scenario like shown in Figure 6 this means that the arrival
38
39 rate varies depending on the PHY mode zone he currently resides in. This behavior is
40
41 illustrated in Figure 3. The amount of stored data in the SCache during the gap period is used
42
43 to “fill up” the otherwise unused resources, reflected by the grey area.
44
45
46

47
48 To model such an arrival process for each PHY mode a state has to be defined. Assigned to
49
50 this state is a data rate λ_i and from the before described mobility simulations transition rates p_{ij}
51
52 from state i to state j can be derived. This perfectly fits to a MMPP arrival process. Such a
53
54 MMPP process can be translated to a MAP as shown in Equation (1).
55
56

$$57$$

$$58$$

$$59$$

$$60$$

$$\Gamma = \begin{pmatrix} -p_1 & p_{12} & \cdots & p_{1N} \\ p_{21} & -p_2 & \cdots & p_{2N} \\ \vdots & \vdots & \ddots & \vdots \\ p_{N1} & p_{N2} & \cdots & -p_N \end{pmatrix}, p_i = \sum_{j \neq i} p_{ij}, \Lambda = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_N \end{pmatrix} \quad (10)$$

1
2
3 But as shown in Figure 14 the data rate might drop to the backbone rate. Two reasons can be
4 cited: either the SCache is drained so that only packets can be transmitted which arrive in the
5 SCache from the video server or the TB is full. Although this does not mean that the data rate
6 decreases the amount of packets accepted from the queue drops to the playback rate and
7 simultaneously the loss ratio of the system increases. Only if a packet leaves the TB a new
8 packet can be stored. Otherwise it has to be discarded. But in such a case the arrival process
9 does not have to be adapted since it has no influence on the idle probability of the queue
10 whether the arrival rate drops down or packets are discarded. Although the loss rate of the
11 queuing system increases, the idle probability remains unchanged.

12
13 So the question is, what happens if the SCache is drained? As shown in Figure 4, the
14 combination of SCache and TB can be seen as a black box system with an incoming data rate
15 λ and an outgoing data rate which is as well equal to λ . This means that after the playback of
16 the video has started the amount of buffered data in the black box does not (or only slightly)
17 change.

18
19 For later analysis now it is assumed that the maximum TB size is equal to the (initial) fill
20 level of this black box (a bigger TB would even improve the performance so that the later
21 results are a lower limit). Under this condition the occurrence of an empty SCache coincides
22 with a full TB. Therefore the aforementioned statement holds that the arrival process can be
23 modeled without any drop of the data rate. The situation that the data rate drops to the BL is
24 implicitly included in the queuing model by discarding packets when the queue is full.

25
26 If the full capacity share is not available for the video streaming process then the data rates in
27 each PHY mode state have to be adapted accordingly. If only 50% of the network resources
28 can be used by the streaming process than the data rate in each PHY mode has to be halved as
29 well.

30
31 To model a backbone limited network, the same MMPP based approach can be used as for the
32 SC-enabled system. It is enough to cut in each state the data rate to the value of the BL. But

1
2
3 the adaptation due to limited network resources for the SC enabled service has to take place
4
5 ahead.
6
7

10 7. Terminal Buffer Size

11
12 With the help of the before described finite buffer queue it is possible to derive the idle
13 probability of the end user device depending on its buffer size. While the analysis with the
14 MMAP/G/1 queuing system has allowed an estimation of the initial buffer fill level now the
15 overall Terminal Buffer size is discussed.
16
17

18
19 In Figure 15 the idle probability depending on the capacity share is shown. The idle
20 probability reflects the likelihood that the TB is drained and the last chunk of video data in the
21 queue is played back. On the x -axis the available capacity share of the wireless network is
22 shown. The parameter of the set of curves is the TB size. The last parameter is measured in
23 seconds the video playback can be continued with a completely filled TB and no additionally
24 arriving packets.
25
26

27
28 It can be observed that the available capacity share has a deep impact on the performance of
29 the SC enhanced network setup. For a realistic service provisioning the idle probability should
30 be below 5%. In that case video playbacks are seldom interrupted and the user has the
31 impression of continuous network connectivity.
32
33

34
35 Only for a very large TB size of 300 seconds, the system remains within reasonable
36 boundaries even if only a capacity share of 25% is taken for the video streaming process. For
37 smaller dimensioned TBs the impact on the idle probability is clearly visible. While for
38 capacity shares of more than 80% all curves are below 7% idle probability, and therefore at
39 least close to the prior defined quality level, the situation changes in case of lower capacity
40 values. The idle probability increases substantially so that the service quality is no longer
41 sufficient. Especially for a TB size of 100 seconds an idle probability of 5% can only be
42 reached if the full capacity share of 100% is reserved for the regarded service.
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 If these results are compared with the outcomes of the *MMAP/G/I* queuing system it can be
4 shown that they match very well. In Table 2, the 95th percentile of the *MMAP/G/I* and the idle
5 probability of the *MAP/G/I/N* systems are compared. In the latter a TB size is assumed which
6 corresponds to the 95th percentile of the first analysis.
7
8
9

10
11 It can be seen that a buffer size which is derived from the 95th percentile produces an idle
12 probability of around 5%. Only for heavy loaded systems (where only one-fourth of the
13 capacity remains for the video streaming) do the results diverge. But the reason is that in the
14 *MMAP/G/I* case it was assumed that during coverage the system always succeeds to drain the
15 SCache. In an overload situation this assumption might be wrong. Gap and coverage period's
16 size differ around an average value so that in case of a long gap and a following short
17 coverage period it might occur that some packets remain in the SCache when the connection
18 breaks again. Then these packets have to wait until the next coverage period before getting
19 delivered. Therefore, the waiting time is extended further and the 95 percentile is pushed to
20 even higher dimensions.
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35

36 The earlier results show a good match between the two analytical approaches. However, a
37 simulative evaluation could strengthen the results. For the considered scenario, very long
38 simulation runs are necessary to cover enough coverage gaps and reach statistically
39 significant results. Together with the sophisticated modeling of the arrival process, the
40 simulation effort would be too complex to reach results within a reasonable amount of time,
41 so the analytical approach has to be taken. Nevertheless for simplistic scenario setups
42 simulations were conducted. A very good match of the packet delay CDFs between the
43 simulation and the analytical results could be shown. As the results are of no further relevance
44 for the here regarded scenario, they are omitted.
45
46
47
48
49
50
51
52
53
54
55
56

57 To dimension a SC network from the analysis of the *MMAP/G/I* queuing system the initial fill
58 level of the TB can be deduced. This guarantees that the video playback does not stop in 95%
59 of the cases during the first gap period. Such a value is already a first indication for the
60

1
2
3 required time shift. Nevertheless, the finite buffer analysis delivers the exact results as it
4 includes not only the first gap but also all subsequent gaps. Thus, the long term probability
5
6
7
8 can be derived so that the video playback does not stop even during longer video streams. It is
9
10 important to note that these two values, at least for low loaded networks, match so that the
11
12 M/G/1/N based analysis is not necessary in each case.
13
14
15
16

17 **8. Motorway Application Scenario**

18
19 To further illustrate the applicability of Smart Caching and its capability to increase
20 performance in delay tolerant networks a second scenario is investigated. In this motorway
21 scenario users go by car and access again a high quality video stream. Due to the massive
22 deployment costs not the complete motorway is covered with broadband wireless access, but
23
24 only in certain intervals the APs are mounted, see Figure 16. The AP sites are chosen
25 depending on the motorway infrastructure. To provide power supply and backbone
26 connectivity usually places like bridges which cross the motorway are the best places. While
27 moving from an AP site to another AP site data is downloaded and stored in the Terminal
28 Buffer when enough network capacity is available – hence, close to the node. This data is then
29 consumed during the idle periods between two AP sites. In such phases the streaming of data
30 is continued between video server and Smart Cache.
31
32
33
34
35
36
37
38
39
40
41
42
43
44

45 In such a scenario setup it can be assumed that the attenuation between sender and receiver is
46 reduced to a minimum. Therefore the attenuation factor γ is set to a free space value of 2. The
47 average distance between the APs is chosen in such a way that the wireless network coverage
48 is equal to 50% in a first setup and reduced to 20% in a second scenario.
49
50
51
52

53 The probability density of the average car velocity is taken from [19] and corresponds to a
54 normal distribution. The density of the gap distance between two concurrent coverage zones
55 is as well normally distributed. All parameters are listed in Table 3.
56
57
58
59
60

1
2
3 From the distributions of the car velocity and the gap distance, the probability density of the
4 duration of a gap period can be derived by applying the following formula.
5
6

$$f_t(t) = \int_0^{\infty} v f_l(vt) f_v(v) dv \quad (11)$$

7
8
9
10
11 In Figure 17 the average data rate with (dot-dashed) and without SC (solid) functionality for
12 such a scenario setup depending on the Backbone Limit is shown (capacity share varied from
13 10 to 100% in 15% steps). On the left, the coverage ratio is 50% and on the right, it is 20%.

14
15
16
17
18 Now it is obvious that especially for low coverage ratios the performance gain of the SC-
19 enabled networks is substantial. For 20% network coverage the required average data rate of
20 3.37 Mbit/s can be provided for an available capacity share of 90 – 100%. With SC the BL
21 does not have to be much above the average value of 3.37 Mbit/s while with a legacy network
22 setup this value dramatically increases to more than 40 Mbit/s. As expected, SC shows its
23 capability for substantial performance enhancement in low coverage scenarios.
24
25
26
27
28
29
30
31

32
33 But before the TB size is discussed an additional improvement of SC is introduced. It is
34 termed “Making a virtue out of necessity”. Delay tolerant networks have to deal with
35 connection interruptions. It has been shown that with SC, connectivity gap periods of more
36 than 100 seconds can be bridged for services like video streaming. But if such long periods
37 have to be handled it might be useful to even lengthen this period. Why should we do so? The
38 outer part of the coverage region of a WiMAX AP is served by a very robust PHY mode to
39 give the receiver a chance despite the long distance to decode the signal. But this also means
40 that the data rate is simultaneously decreased. The transfer of 1 MB of data requires in this
41 robust PHY mode more than 1.3 seconds while close to the AP this can be completed in less
42 than 0.15 seconds. So from a resource allocation point of view it is useful to wait until a
43 mobile user reaches the most inner PHY mode region before massive data transfers take
44 place.
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

Such a behavior can be obtained by deactivating the most robust PHY modes which cover the
outer annuli of a coverage region. The more PHY modes are deactivated the longer the gap

1
2
3 between concurrent AP zones becomes. But simultaneously the efficiency of the network
4
5 increases.
6

7
8 In Figure 18 the CDF of the packet waiting time is shown depending on how much PHY
9
10 modes are deactivated. It can be seen that the 95th percentile which reflects the initial buffer
11
12 fill level – but is as well a good approximation of the required TB size – increases the more
13
14 PHY modes are deactivated. But the increase is only moderate from 220 to 380 seconds. The
15
16 impact this has on the network efficiency can be seen on Figure 18 (right). For different
17
18 values of traffic intensity the ratio of possible video users per 100 cars is shown depending on
19
20 the number of deactivated PHY modes. The reason for the efficiency gain is not that more
21
22 users can be served by one AP, but that less users have to download data simultaneously. If
23
24 only the inner PHY mode zones are used for the transfer of data then the amount of cars
25
26 residing in the coverage area decreases. This means that the ratio of supported users can be
27
28 increased accordingly.
29
30
31
32

33
34 For such a setup there has to be found a tradeoff between the packet waiting time, its resulting
35
36 initial buffer fill level and the network performance gain. But of course it is possible to switch
37
38 between the different configurations. At the beginning a mobile user starts with a fill level of
39
40 220 seconds and all PHY modes are used. As soon as the fill level of the buffer could be
41
42 increased above 280 seconds (e.g., due to very good network performance) the most robust
43
44 PHY mode is deactivated. This can continue until 5 PHY modes are switched off. Of course
45
46 the reverse is also possible. If there is a severe threat that the TB drains then the outer PHY
47
48 modes can be reactivated. The exclusive usage of only the inner PHY mode (mode 1-6
49
50 deactivated) is not possible as the network resources are insufficient to provide the necessary
51
52 average data rate in case of such a reduced coverage ratio.
53
54
55
56
57
58
59
60

9. Summary and Conclusions

Smart Caching is introduced as a method for making networks delay tolerant and handling connectivity disruptions. The technique allows the optimization of throughput capacity in broadband wireless networks. Moreover, it supports and elaborates the extensive buffering of service data in the end user device. This allows the virtual continuation of broadband services even with intermittent network connectivity.

The complete communication path between server and client is included so that it is possible to assess the potential of Smart Caching for an end-to-end service provisioning. In the analysis, a video streaming service is regarded as it has the highest demands to communication networks. Nevertheless, the content to be cached is easily predictable so that taking into account deviations is not required. For further evaluations this constraint might be excluded as other services do not allow for a precise prediction.

The presented results allow a good dimensioning of streaming services as the delay analysis gives the initial time shift, and the analysis of the buffer size allows the dimensioning of the Terminal Buffer in the user terminal. For services that are delay tolerant by default (such as file transfer), the delay evaluations are not of much interest but the achievable data rate in Figure 8 and Figure 17 can be perfectly used to derive the necessary transfer durations. So this article covers a wide range of services. In the future, it is necessary to investigate services which are not 100% predictable.

Two application scenarios are presented – a pedestrian user in an urban environment and a vehicular user on a motorway. In both scenarios the applicability of Smart Caching is proven. The analysis with the support of two sophisticated queuing models allows a very accurate dimensioning of the Terminal Buffer size. Furthermore the requirements of an initial fill level for the Terminal Buffer is derived and verified by the analysis.

In the motorway scenario a further improvement of Smart Caching is given. It allows an enhancement of the network capacity so that more users can be supported simultaneously.

1
2
3 The additional delay which arises can be absorbed by the Smart Caching technology. All this
4
5 proves that Smart Caching is a perfect way to circumvent problems occurring in wireless
6
7 networks with patchy coverage.
8
9
10
11
12
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

For Peer Review

References

- [1] H. Balakrishnan, S. Seshan, E. Amir, and R.H. Katz, *Improving TCI/IP performance over wireless networks*, In Proc. 1st ACM Int'l Conf. on Mobile Computing and Networking (Mobicom), pp. 2—11 , Nov. 1995.
- [2] M. Zink, A. Jonas, C. Gridwodz, and R. Steinmetz, *LC-RTP (loss collection RTP): reliability for video caching in the Internet*, In Proc. of Seventh International Conference on Parallel and Distributed Systems, pp. 281—286 , 2000.
- [3] L. Song, D. Kotz, R. Jain, and X. He, *Evaluating Next-Cell Predictors with Extensive Wi-Fi Mobility Data*, In IEEE Transactions on Mobile Computing, pp. 1633—1649, 2006.
- [4] S. Podlipnig and L. Böszörmenyi, *A survey of Web cache replacement strategies*, In ACM Computing Surveys (CSUR), Volume 35, Number 4, pp. 374—398 , 2003.
- [5] Z. Su, Q. Yang, and H.J. Zhang , *A prediction system for multimedia pre-fetching in Internet*, In Proceedings of the eighth ACM international conference on Multimedia, pp. 3—11 , 2000.
- [6] K. Fall and S. Farrell, *DTN: an architectural retrospective*, In IEEE Journal on Selected Areas in Communications, Volume 26, Number 5, pp. 828—836 , 2008.
- [7] M. Demmer, E. Brewer, K. Fall, S. Jain, M. Ho, and R. Patra, *Implementing Delay-Tolerant Networking*, Intel Research, Berkeley, Technical Report, IRB-TR-04-020, Dec. 2004.
- [8] Y. Lin, B. Li, and B. Liang, *Stochastic analysis of network coding in epidemic routing*, In IEEE Journal on Selected Areas in Communications, Volume 26, Number 5, pp. 794—808 , 2008.
- [9] X. Zhang, J.K. Kurose, B.N. Levine, D. Towsley, and H. Zhang, *Study of a bus-based disruption-tolerant network: mobility modeling and impact on routing*, In Proceedings of the 13th annual ACM international conference on Mobile computing and networking, pp. 195—206 , 2007.
- [10] R.A. Nichols and A.R. Hammons, *Performance of DTN-Based Free-Space Optical Networks with Mobility*, In Proceedings of IEEE Military Communications Conference 2007, MILCOM 2007, pp. 1—6 , 2007.
- [11] B. Walke, *On the importance of WLANs for 3G cellular radio to become a success*, In Proc. 10th Aachen Symposium on Signal Theory - Mobile Communications, 2001.
- [12] C. Hoymann, *Analysis and performance evaluation of the OFDM-based metropolitan area network IEEE 802.16*, Computer Networks, 2005.
- [13] S. Asmussen, O. Nerman, and M. Olsson, *Fitting phase-type distributions via the EM algorithm*, *Scandinavian Journal of Statistics*, 1996, 23, 419-441.
- [14] P. Seeling, F. Fitzek, and M. Reisslein, *Video Traces for Network Performance Evaluation*, Springer, 2007.
- [15] T. Takine, *Queue Length Distribution in a FIFO Single-Server Queue with Multiple Arrival Streams Having Different Service Time Distributions*, *Queueing Systems*, Vol. 39, No. 4, pp. 349—375, 2001.
- [16] T. Takine, *The Nonpreemptive Priority MAP/G/1 Queue*, *Operations Research*, Vol. 47, No. 6, 1999, pp. 917-927.
- [17] B. Maglaris, D. Anastassiou, P. Sen, G. Karlsson, and J.D. Robbins, *Performance models of statistical multiplexing in packet videocommunications*, *IEEE Transactions on Communications*, Vol. 36, No. 7, pp. 834—844, 1988.
- [18] U. Gupta and P. Laxmi, *Analysis of the MAP/G(a, b)/1/N Queue*, *Queueing Systems*, Springer, 2001, 38, 109-124.
- [19] C.-H. Rokitanski, Editor, *Validation of the Mobility Model of SIMCO2 with Dutch Motorway Measurements*, Prometheus Workshop on Simulation, December 1992.

1
2
3
4 **Proof of Equation (5)**
5
6
7

8 $pdf(t) = P(W = t)$: Probability that the additional delay of an arbitrary packet is t
9

10 $p(x) = P(gap_duration = x)$: Probability that gap duration is x
11

12 $P(Packet_in_gap)$: Probability that the arrival of the packet occurs when the user is in a
13 gap
14
15

16
17 Three factors have to occur in order to perceive an additional delay t . First there must be a
18 connectivity gap of size x . Second the arrival of the packet has to occur while the terminal is
19 in the gap. And finally the condition has to hold that the gap duration x must be larger than the
20 delay t .
21

22
23 $P(W = t) = P(Packet_in_gap, gap = x | t \leq x)$
24

25 Since the condition $t \leq x$ is independent of the other probabilities it can be stated that
26

27
28 $P(W = t) = P(Packet_in_gap, gap = x, t \leq x)$
29
30 $= \int_y^{\infty} P(Packet_in_gap, gap = x, t = x) dx$
31
32 $= \int_{x'}^{\infty} P(t = x | Packet_in_gap, gap = x) P(Packet_in_gap | gap = x) P(gap = x) dx$
33
34
35

36 Two terms have to be further evaluated $P(Packet_in_a_gap | gap = x)$ and
37 $P(t = x | Packet_in_gap, gap = x)$. The second one describes the probability that if the packet gets
38 definitely an additional delay due to a gap of duration x that the delay t is equal to x . Since it
39 can be assumed that packets arrive equally distributed over the whole gap interval x the
40 probability density is as well equally distributed with:
41
42

43 $P(t = x | Packet_in_gap, gap = x) = \frac{1}{x}$.
44
45

46
47 The other required term $P(Packet_in_a_gap | gap = x)$ describe the probability density that a
48 packet arrives during the gap given the fact the size of the gap is x . This cannot be directly
49 evaluated. But with the same reasoning as before it is clear that the probability that the packet
50 arrives during the gap increases linearly with the gap duration x .
51

52 $P(Packet_in_gap | gap = x) = cx$
53
54

55 where c is just a constant factor. In order to determine c it can be used that the integral of the
56 unconditional probability $P(Packet_in_a_gap, gap = x)$ must be equal to 1 .
57

58 $\int_0^{\infty} P(Packet_in_gap, gap = x) dx = \int_0^{\infty} P(Packet_in_gap | gap = x) P(gap = x) dx$
59
60 $= \int_0^{\infty} cx P(gap = x) dx = c \int_0^{\infty} x P(gap = x) dx = cE[x] = 1$

The last transformation is just the expected value of the gap duration probability distribution. So it follows that:

$$P(\text{Packet_in_gap} | \text{gap} = x) = \frac{1}{E[x]} x.$$

With all this $P(W=t)$ results to

$$\begin{aligned} P(W=t) &= \int_{x'}^{\infty} P(t=x | \text{Packet_in_gap}, \text{gap} = x) P(\text{Packet_in_gap} | \text{gap} = x) P(\text{gap} = x) dx \\ &= \int_{x'}^{\infty} \frac{1}{x} \frac{x}{E[x]} p(x) dx = \int_{x'}^{\infty} \frac{p(x)}{E[x]} dx \text{ q.e.d.} \end{aligned}$$

1
2
3
4
5
6
7
8
9
10
11
12
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

	C(0)	a	Data Rate [$10^6/s$]	M
Video (MPEG4)	0.509	0.09	3.37	20

Table 1

For Peer Review

Net ratio/buffer size (based on 95 percentile) [s]	Idle probability
1.0 / 102	4.37%
0.5 / 112	6.92%
0.25 / 142	11.78%

Table 2

For Peer Review

1
2
3
4
5
6
7
8
9
10
11
12
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

Car velocity	Gap size (Scenario type 1)	Gap size (Scenario type 2)
$f_v(v)$	$f_l(l)$	$f_l(l)$
Normally distributed	Coverage=50%	Coverage=20%
$\mu=104$ km/h	Normally distributed	Normally distributed
$\sigma=12.5$ km/h	$\mu=5100$ m	$\mu=20.5$ km
	$\sigma=500$ m	$\sigma=5$ km

Table 3

For Peer Review

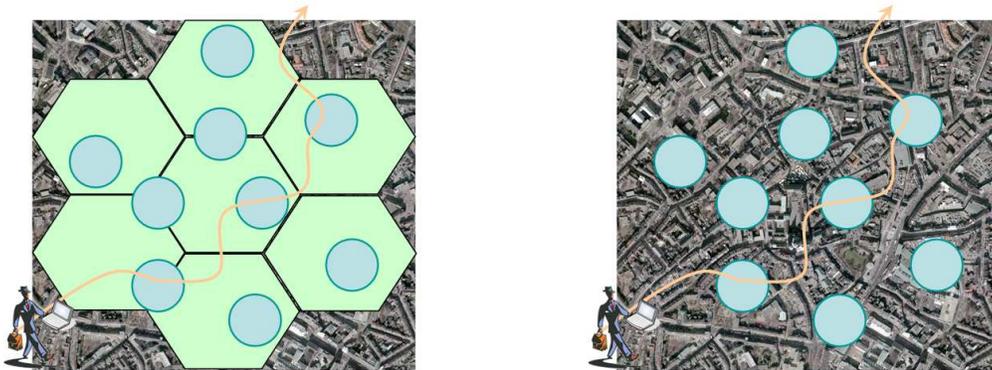


Figure 1: Coverage of heterogeneous wireless networks
(Green hexagons represent UMTS coverage and blue circles WiMAX Hot Spots)

213x86mm (150 x 150 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
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

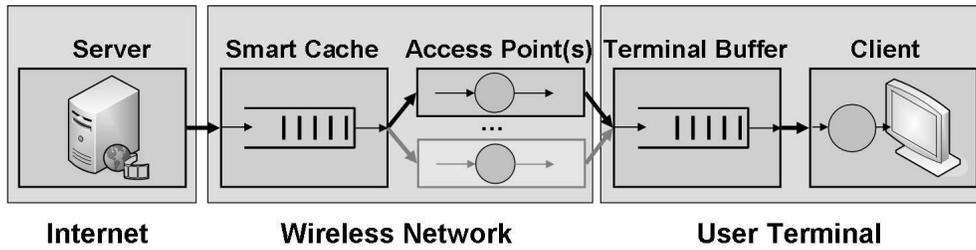


Figure 2: Architecture of Smart Caching
253x71mm (150 x 150 DPI)

Or Peer Review

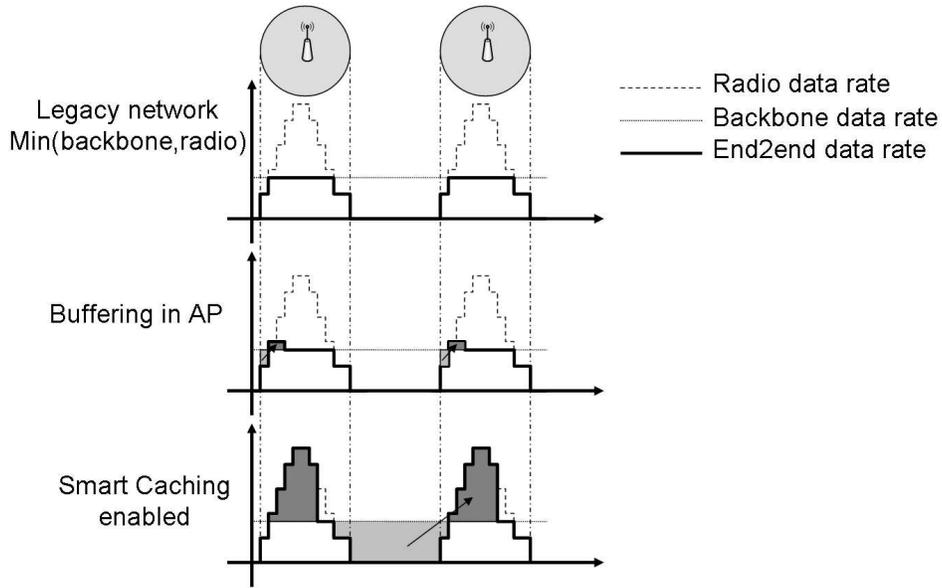


Figure 3: Impact of Smart Caching on Data Rate
243x154mm (150 x 150 DPI)

Review

1
2
3
4
5
6
7
8
9
10
11
12
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

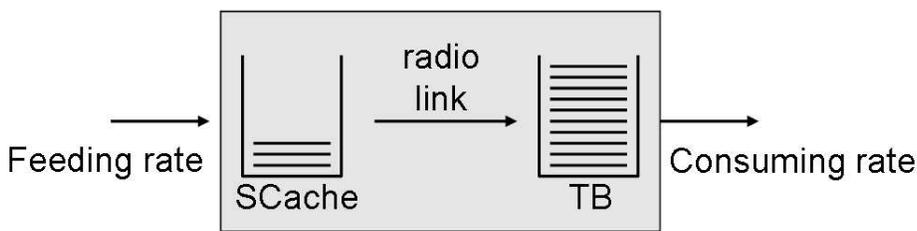


Figure 4: Buffering Concept of Smart Caching
181x47mm (150 x 150 DPI)

For Peer Review

1
2
3
4
5
6
7
8
9
10
11
12
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

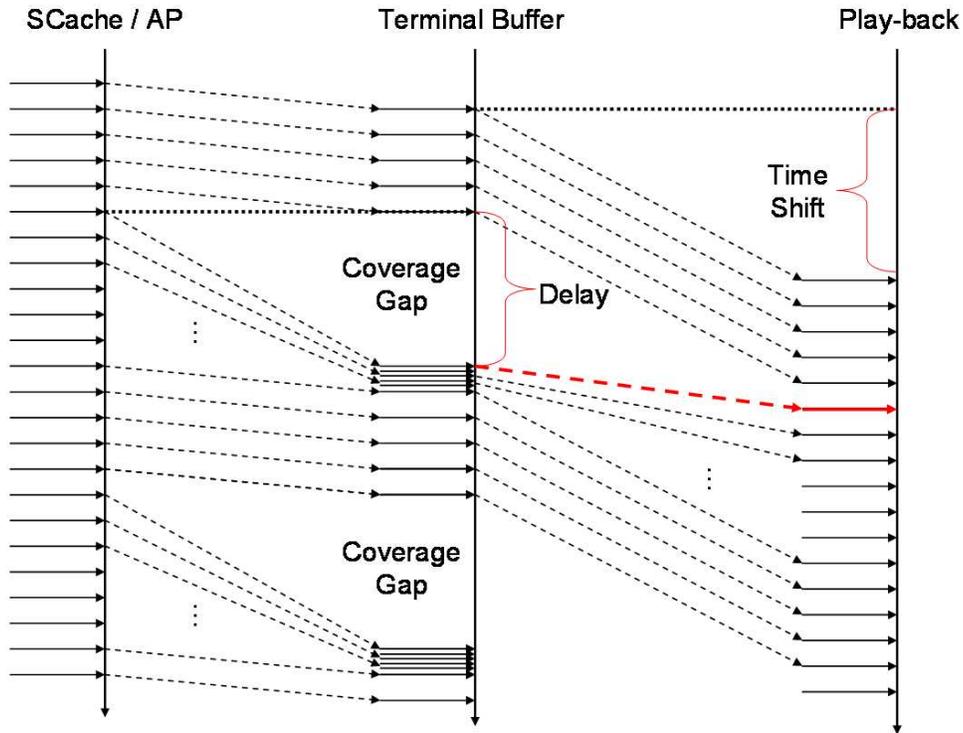


Figure 5: Relation between Time Shift and Packet Delay
185x140mm (150 x 150 DPI)

Review

1
2
3
4
5
6
7
8
9
10
11
12
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

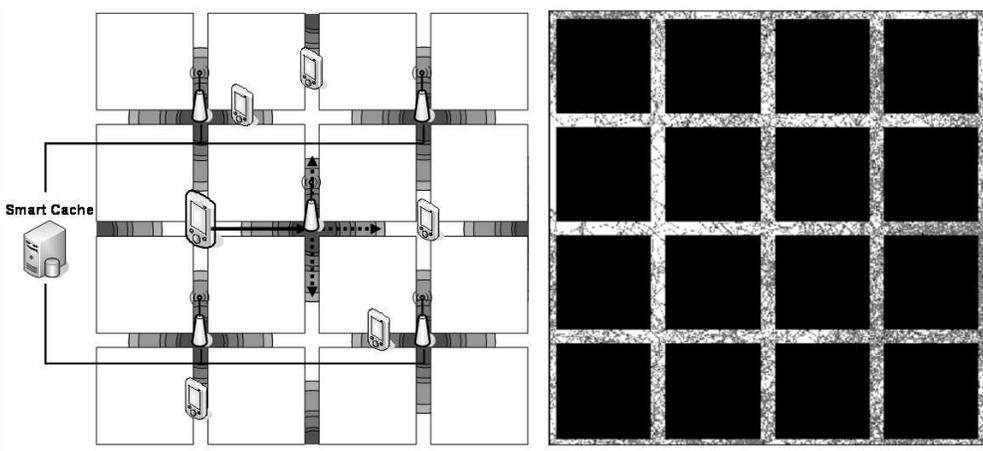


Figure 6: Urban Scenario and Mobility Simulation
197x90mm (150 x 150 DPI)

Peer Review

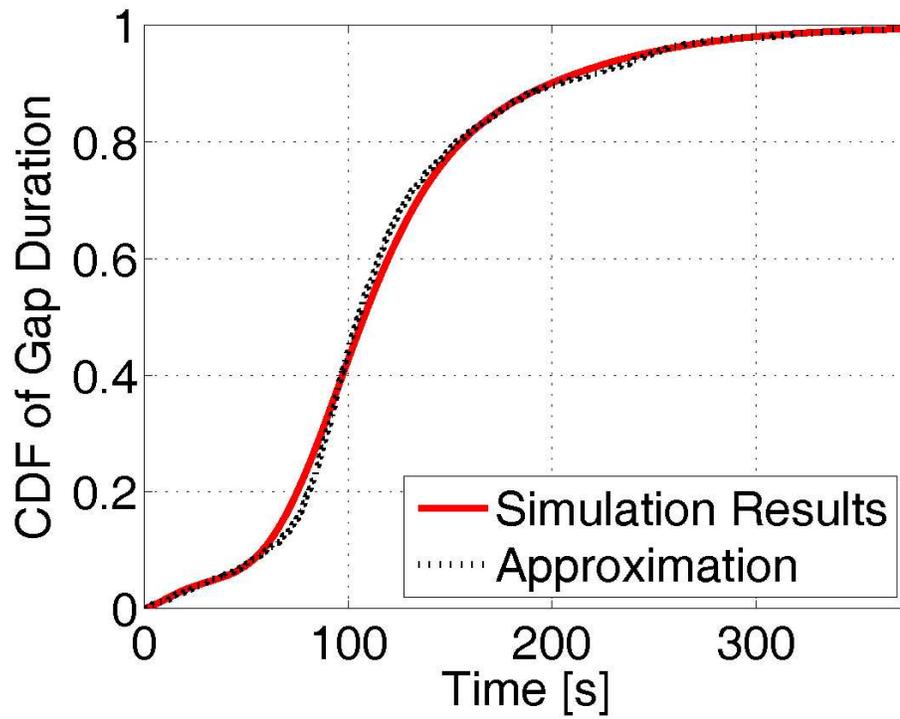


Figure 7: Gap Duration Distribution
148x111mm (200 x 200 DPI)

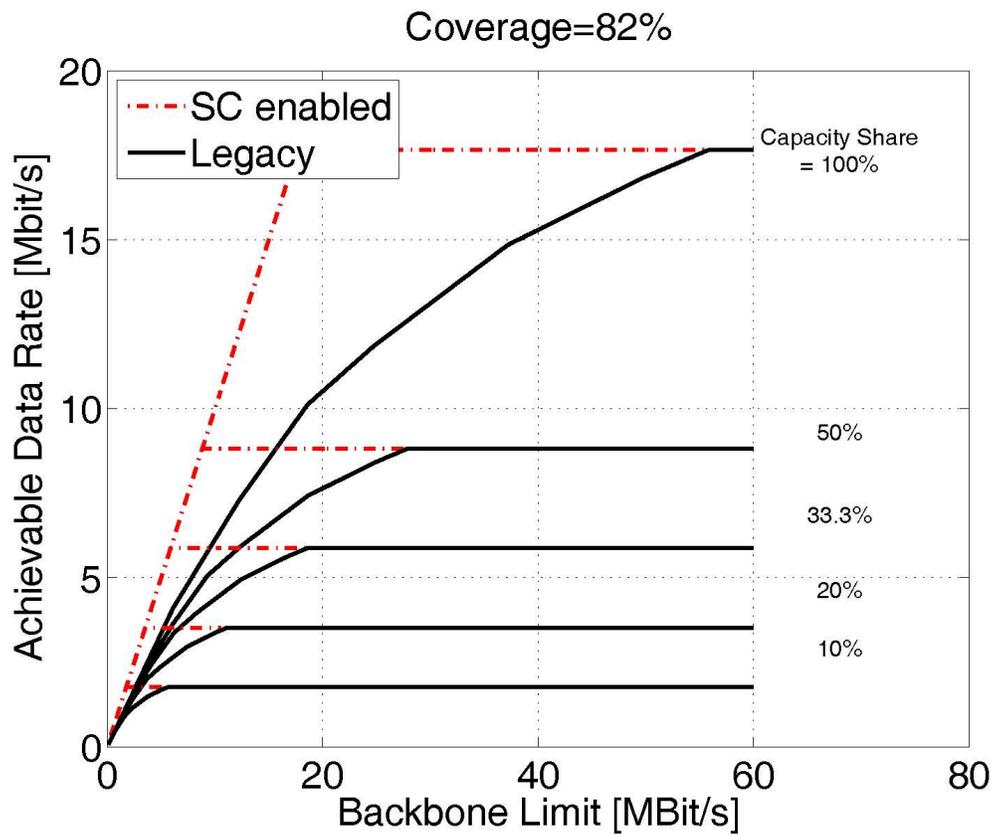


Figure 8: Data Rate depending on Backbone Limit
182x153mm (200 x 200 DPI)

view

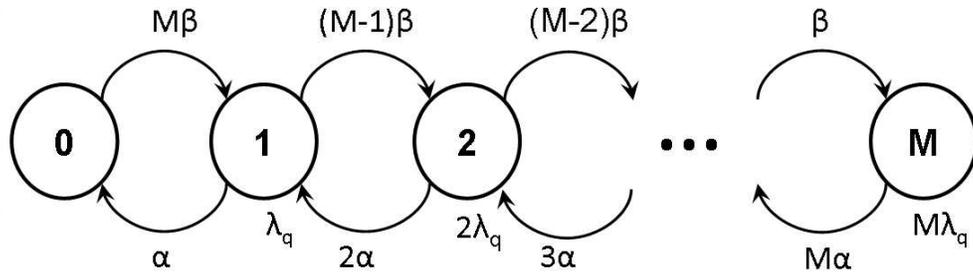


Figure 9: MMPP model for Video Streaming
188x70mm (150 x 150 DPI)

Peer Review

1
2
3
4
5
6
7
8
9
10
11
12
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

1
2
3
4
5
6
7
8
9
10
11
12
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

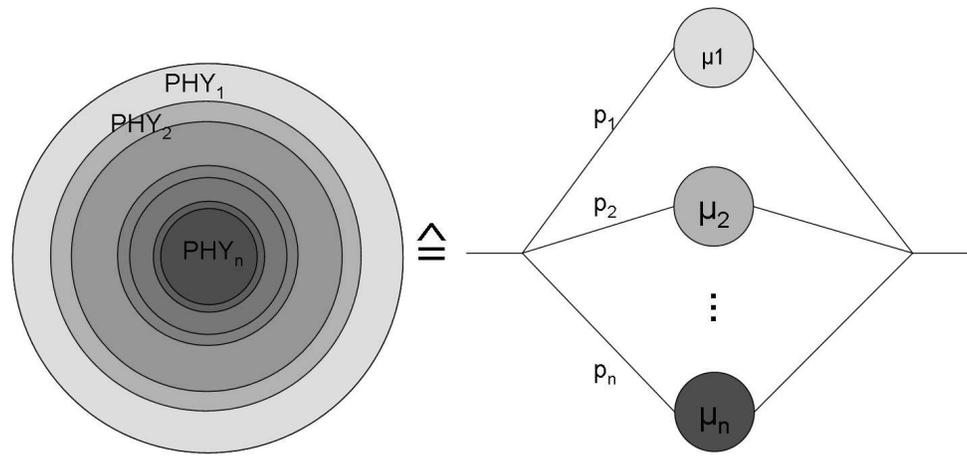


Figure 10: Service Process Modeling
250x123mm (150 x 150 DPI)

Peer Review

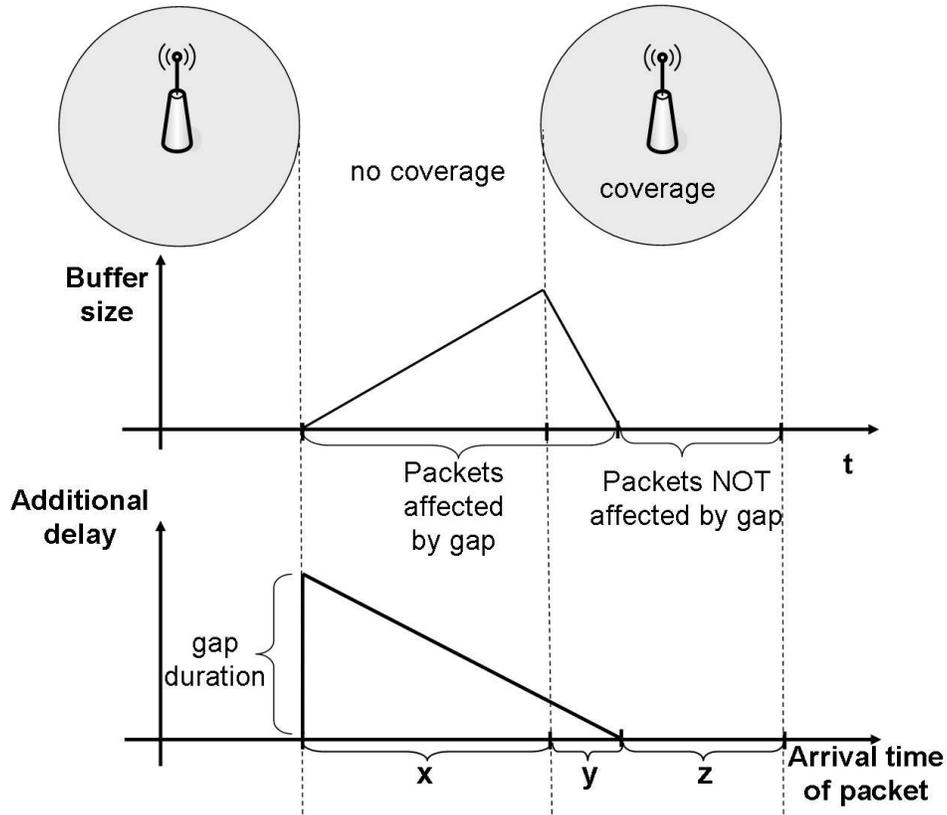


Figure 11: Influence of Coverage Gap on Buffer Size and Packet Delay
 224x191mm (150 x 150 DPI)

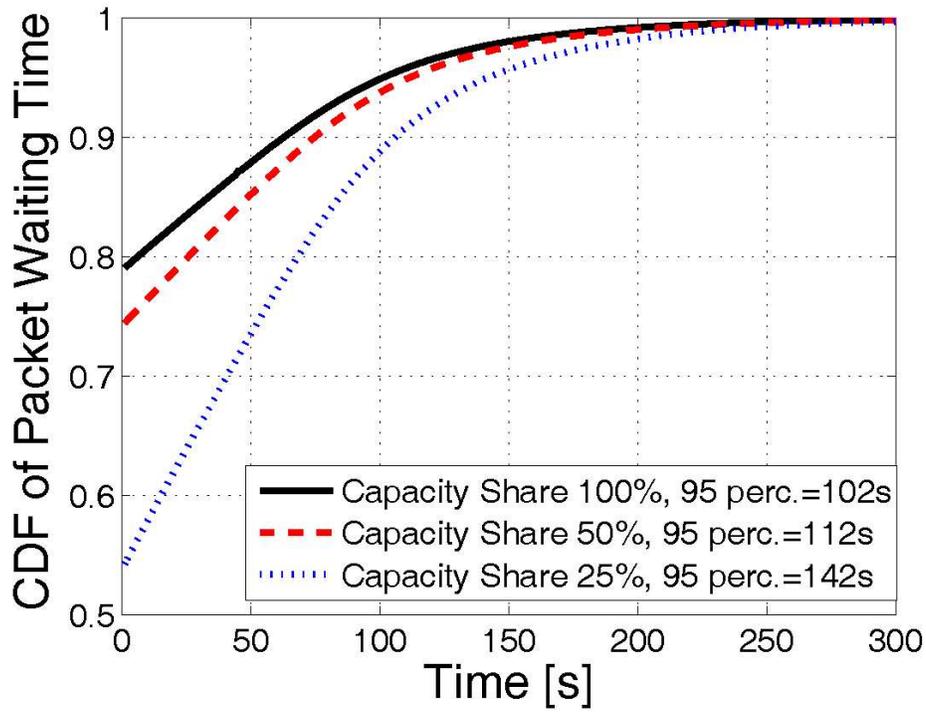


Figure 12: Packet Waiting Time in Smart Caching-enabled Urban Scenario
148x111mm (200 x 200 DPI)

Review

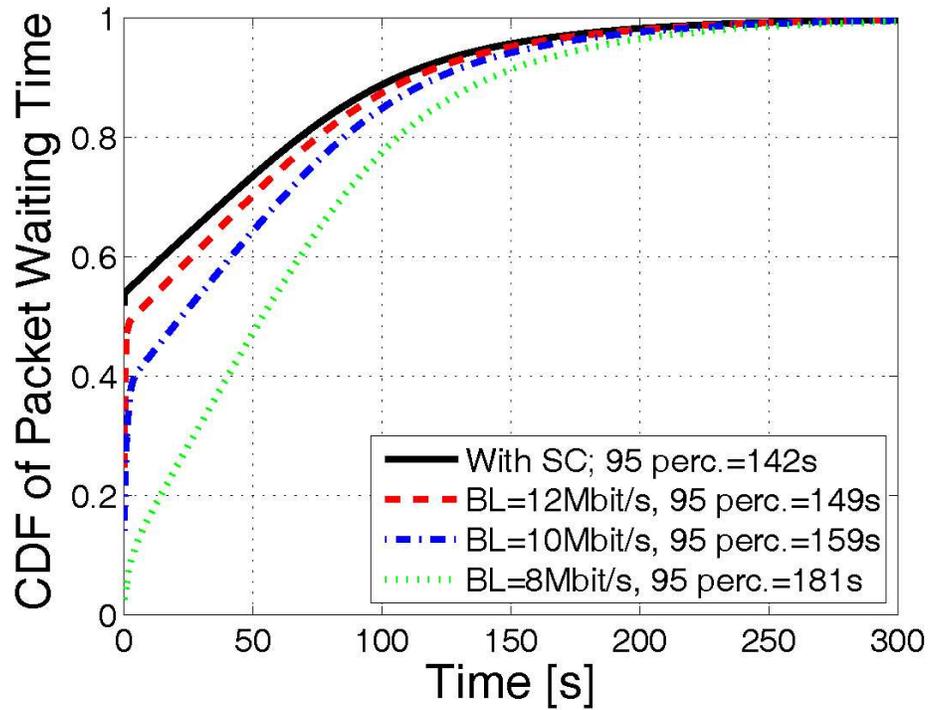


Figure 13: Impact of BL on CDF
148x111mm (200 x 200 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
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

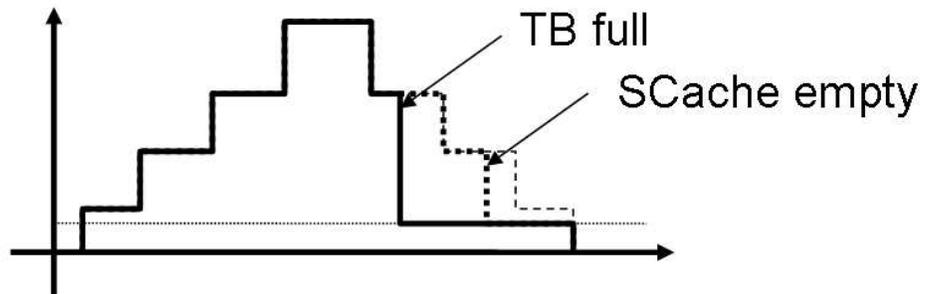


Figure 14: Impact of full TB and empty SCache
138x50mm (150 x 150 DPI)

Peer Review

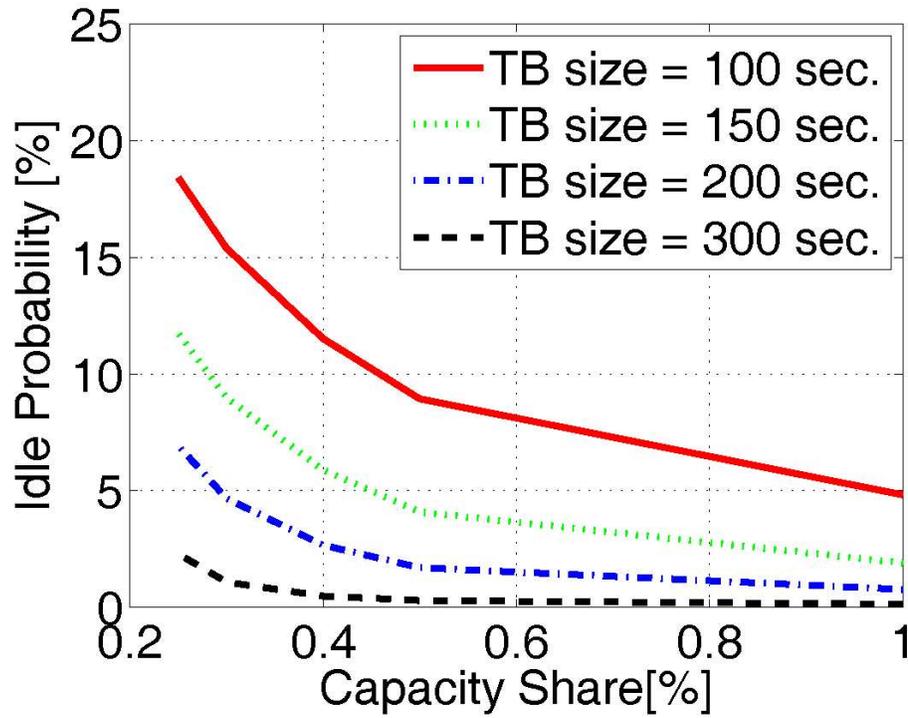


Figure 15: Impact of TB Size on Idle Probability
148x111mm (200 x 200 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
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

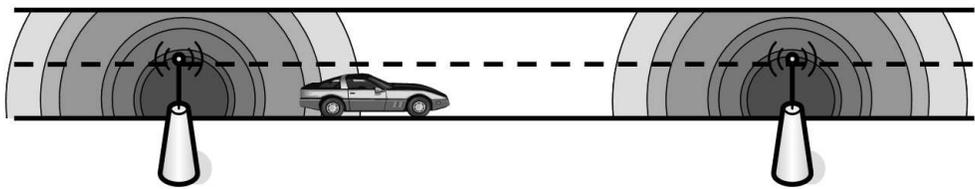


Figure 16: Highway Scenario
219x47mm (150 x 150 DPI)

For Peer Review

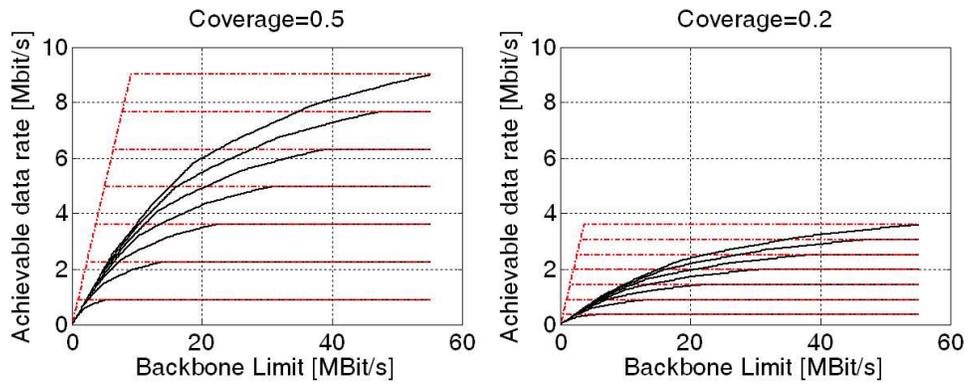


Figure 17: Data Rate depending on Network Coverage
209x83mm (150 x 150 DPI)

1
2
3
4
5
6
7
8
9
10
11
12
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

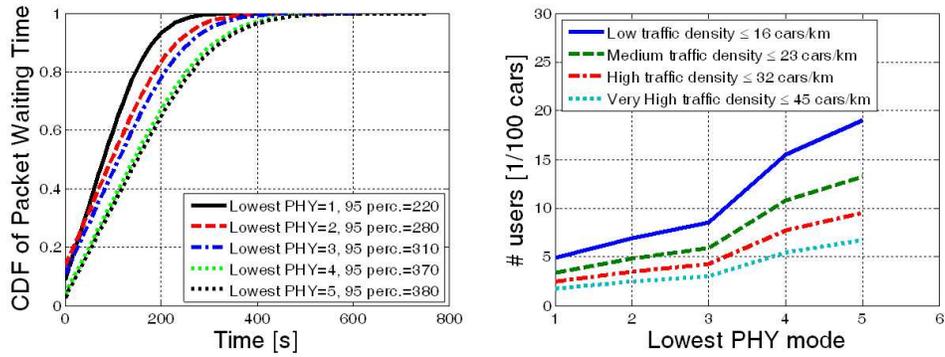


Figure 18: Capacity Gain by PHY Mode Deactivation
214x79mm (150 x 150 DPI)