

(19)



(11)

EP 3 451 628 A1

(12)

EUROPEAN PATENT APPLICATION

(43) Date of publication:
06.03.2019 Bulletin 2019/10

(51) Int Cl.:
H04L 29/08 (2006.01) H04B 7/0413 (2017.01)
H04B 7/024 (2017.01)

(21) Application number: **17290111.8**

(22) Date of filing: **31.08.2017**

(84) Designated Contracting States:
AL AT BE BG CH CY CZ DE DK EE ES FI FR GB GR HR HU IE IS IT LI LT LU LV MC MK MT NL NO PL PT RO RS SE SI SK SM TR
Designated Extension States:
BA ME
Designated Validation States:
MA MD

(72) Inventors:
• **Elia, Petros**
06160 Juan Les Pins (FR)
• **Lampiris, Eleftherios**
06410 Biot (FR)

(74) Representative: **Schuffenecker, Thierry**
120 Chemin de la Maure
06800 Cagnes sur Mer (FR)

(71) Applicant: **Institut Eurecom G.I.E.**
06410 Biot-Sophia Antipolis (FR)

(54) **SYSTEM AND METHOD FOR MANAGING DISTRIBUTION OF INFORMATION IN MULTI-ANTENNA AND MULTI-TRANSMITTER ENVIRONMENTS**

(57) This patent relates to a communication between a wireless base station with L transmit antennas and a set of receiving nodes equipped with caches so that content from a set of files can be stored in these caches. By formulating a novel way of setting users in groups and a way to store content partially as function of this grouping, this patent allows for increased gains in terms of number of users served simultaneously compared to the state-of-art. The approach taken here makes use of the multiple transmit antennas to boost the number of multicasting users, in some settings, by up to L times compared to the state-of-art, thus reducing significantly the transmission time. The invented solution allows for caching of file segments in such a way such that users belonging in the same group have high correlation of content, while lower correlation appears between the content inter-group caches.

Moreover, the invention is, also, presented for the settings of i) wired communications, ii) cooperative communications between two or more base stations transmitting to a set of receiving nodes in either a wireless or a wired network and iii) in a device-to-device setting where users desire to exchange files between one another

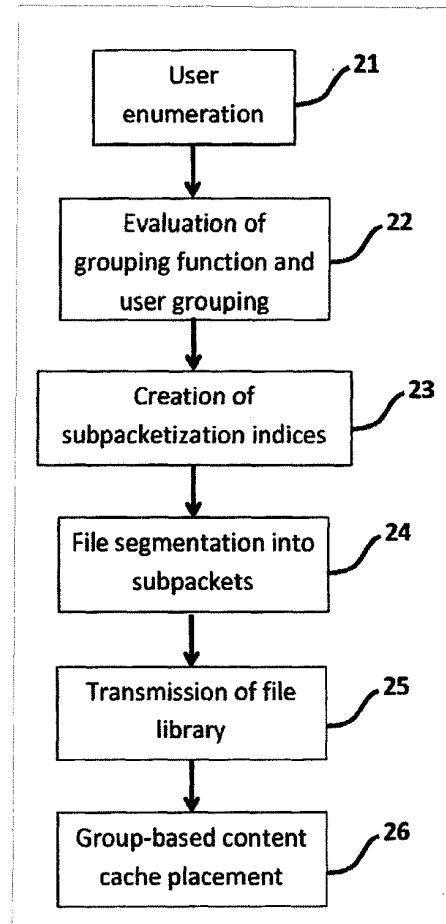


Fig. 2

EP 3 451 628 A1

Description**Technical Field**

5 **[0001]** The present invention relates to the field of digital communications and specifically to a system and a method for managing distribution of information in multi-antenna and multi-transmitter environments.

Background Art

10 **[0002]** Nowadays, extensive research is performed in telecommunications in order to improve the speed and reduce the delay of delivery of files, such as video and audio files. In wireless communications, for instance, a well-known method for improving communication between a base station (B) and different receiving nodes N_1, N_2, \dots, N_K , is the use of multiple antennas at the transmitting base station which essentially employ channel information in order to beamform signals to different users, thus generally allowing each user to receive only their own messages. This ability to send vectors of signals (each signal departing from one of the transmit antennas) at a time, results in the so called 'multiplexing gain', which refers to the ability to serve multiple users at a time. Saying that a multi-antenna (MIMO) precoding system offers multiplexing gain G_M generally implies that the system can serve - over a single time-frequency slot - G_M users at a time. In practice though, the difficulty of obtaining proper channel state information, severely limits the standalone effectiveness of such multi-antenna precoding techniques, see reference [1] below.

20 **[0003]** In wireless, as well as wired communications, another well-known method used to reduce delay in communication between a transmitter and different receiving users, is referred to as "cache-aided multicasting", as discussed in reference [2] below, and it involves the use of storage (also referred to as "caches") in the users' devices. Caches have historically been used to pre-store typically popular content, such as for example, popular movies. Figure 1 illustrates such a known architecture based on a base-station **11** (B) that communicates with nodes **14** through multiple antennas **12** in particular, and more general to nodes N_1, N_2, \dots, N_K , in a wireless communication system, where nodes N_1, N_2, \dots, N_K , have storage capabilities, as illustrated by cache **15** in figure 1. This recent approach, of cache-aided multicasting, works by carefully filling up each receiver cache with content that is possibly undesired to that user but which is expected to be desired by different users, such that when this undesired content is transmitted, the user can use its cache to remove this undesired interference. Receiver-side caches are used in such a way so that the transmitter can simultaneously serve different files to different users with a *composite signal* which carries a combination of what each receiver has requested, and which is multicasted to all receivers simultaneously, such that each receiver uses their own cache in order to remove unwanted files from the composite signal, and thus retrieve their own signal of interest. This approach - in theory - is shown to provide substantial speed-up gains, referred to as "cache-aided multicasting gains". Saying that a multicasting system offers cache-aided multicasting gain G_C generally implies that the system can serve - over a single time-frequency slot - G_C users at a time more than if it had no caches. It was believed that in theory, such cache-aided multicasting systems gains would be very large, inducing a delay that was theoretically calculated to be small, regardless of the number of users requesting files.

30 Nevertheless, soon came the realization that these gains are severely limited by the fact that each file **13** (in figure 1) would have to be segmented into an exponentially large number of subfiles or segments, which severely limits the gains [3], which means that it severely reduces the number of users that can be served at a time, and thus substantially increases communication delay compared to the theoretical result.

Combining multi-antenna precoding with cache-aided multicasting.

45 **[0004]** Given the aforementioned limitations of both multi-antenna precoding and cache-aided multicasting methods, some authors analyzed the possible combination of these two approaches. In theory, this was shown to be possible: combining an L -antenna precoder (which gives multiplexing gain of $G_M = L$) with a cache-aided multicasting method which gives a multicasting gain G_C , it was shown (see [4],[5] below) that one could achieve a theoretical total gain $G_T = L + G_C$. This was interesting because the two methods are seemingly orthogonal: precoding is based on signal separation where each user is expected to get a different signal, while multicasting is based on signal mixing where the users get a common composite signal.

50 **[0005]** However, the authors rapidly identified a main, drastic problem of combining multi-antenna precoding with cache-aided multicasting. To date though, all existing algorithms for combining these two ingredients (MIMO and cache-aided coded multicasting) suffer from an exponentially aggravated astronomical complexity, referred to as sub-packetization or segmentation complexity. This is now a well-known limitation and it has been documented clearly in references [2][4][5] below, and it is now known that this exponential complexity causes severe reductions in the "compromised" cache-aided multicasting gain $\bar{G}_C \ll G_C$ which - under a variety of real-life limitations on the file size, and on

the packet size - is unfortunately not expected to be substantial.

Specifically in a system like our system here, with an L -antenna base station serving K receiving nodes, for this theoretical gain $G_T = L + G_C$ to be achieved, each file would have to be segmented into different parts.

5

$$\binom{K}{G_C} \times \binom{K - G_C - 1}{L - 1}$$

10 **[0006]** One will rapidly observe that, with illustrative values such as:

$K=1000$ users; $G_C = 10$ and $L = 5$ for instance, the formulation above leads to an astronomic figure of 10^{33} , thus ruining any hopes of a realistic, pragmatic solution.

15 **[0007]** Indeed, given that the smallest possible part is the bit, the above expression reveals that each file would need to be of astronomical size for the gains to appear.

[0008] This clearly shows that it is practically infeasible to meaningfully combine multi-antenna precoding methods with cache-aided multicasting methods, because of this aforementioned astronomical segmentation complexity identified in the existing state-of-art algorithms.

20 **[0009]** As a result, *it has been widely accepted in the literature that merging multi-antenna precoding with cache-aided multicasting cannot be implemented in practice* because of this astronomical complexity which *allows only for minimal gains associated to caching*, resulting in unrealistically large communication delays when delivering files of realistic sizes.

[0010] The present invention aims to contradict such theoretical limitation and provides a technically feasible and realistic solution to such problem of combining cache aided multicasting in a MIMO environment.

25 **[0011]** The following Bibliographical references are of interest to this technical problem.

[1] L. Zheng and D. N. C. Tse, "Communication on the Grassmann manifold: A geometric approach to the noncoherent multiple-antenna channel," IEEE Trans. Information Theory, Feb. 2002.

30 [2] M. Maddah-Ali and U. Niesen, "Fundamental limits of caching," IEEE Trans. Information Theory, May 2014.

[3] K. Shanmugam, M. Ji, A. M. Tulino, J. Llorca, and A. G. Dimakis, "Finite length analysis of caching-aided coded multicasting," ArXiv e-prints, Aug. 2015.

35 [4] S. P. Shariatpanahi, A. S. Motahari, and B. H. Khalaj, "Multi-server coded caching," ArXiv e-prints, Aug. 2015.

[5] N. Naderializadeh, M. A. Maddah-Ali, and A. S. Avestimehr, "Fundamental limits of cache-aided interference management," CoRR, vol. abs/1602.04207, 2016. [Online]. Available: <http://arxiv.org/abs/1602.04207>

40 **Summary of the invention**

[0012] It is an object of this invention to combine MIMO and cache-aided multicasting in a way that allows for substantial multicasting gains, with a system that is still technically feasible.

45 **[0013]** It is another object of this invention to introduce a method of using multiple transmit antennas to reduce the complexity associated with caching and cache-aided multicasting algorithms.

[0014] It is still another object of this invention to introduce a method of using multiple transmit antennas to deliver files of limited sizes, with reduced transmission delay.

[0015] It is a further object of this invention to combine caching at the receivers and many antennas at the transmitter, to increase the goodput of modern telecommunication systems.

50 **[0016]** It is still another object of this invention to use caching at the receivers to reduce implementation and hardware complexity of multi-antenna systems.

[0017] It is another object of this invention to introduce a method that extends single antenna caching algorithms to be applied in systems with at least two antennas.

55 **[0018]** It is another object of this invention to introduce a method of using multiple transmit antennas to deliver files, where said method achieves reduced file subpacketization.

[0019] It is another object of this invention to introduce a method of using multiple transmit antennas to deliver files, where said method allows increased packet sizes.

[0020] It is another object of this invention to reduce the complexity associated with decentralized caching and decentralized cache-aided multicasting algorithms.

[0021] It is another object of this invention to increase the gains from using decentralized caching and decentralized cache-aided multicasting algorithms to deliver files of limited size.

5 [0022] It is another object of this invention to apply multiple transmitters rather than just multiple antennas, to achieve the above said objects of this invention.

[0023] It is another object of this invention to apply multiple transmitters or relays to achieve the above said objects of this invention in the presence of a wired network.

10 [0024] The above, and other objects of this invention are achieved by means of a process that applies to a multi-antenna broadcast and multicast system, where a multi-antenna transmitting base station (B) transmits to receivers N_1, N_2, \dots, N_K which are equipped with caches. In this system, the process involves the step of:

- grouping K users into groups containing a specified number of users per group, and
- 15 - involving a first phase or time slot during which the base station employs two or more antennas to transmit information that is received and then partially cached at some or all of the caches of the receivers, the method being a function of at least the grouping of the users.
- involves a second time slot during which some or all receiving nodes proceed to request additional information and based on this as well as on additional channel state information that the receivers feedback to the transmitter, the next function takes place;
- 20 - generating by the base station (B), on the basis of the user grouping, a multi-antenna transmit signal with the purpose of transmitting the requested data files to said receiving nodes;
- decoding by said participating receiving nodes the signal transmitted by said base station (B), using the information received in said first and second time slots and the information in the caches, so as to allow each receiving node N_k to compute an estimation of each requested file F_{rk} .

25 [0025] This system combines multi-antenna precoding with cache-aided multicasting, to reduce the complexity associated with cache-aided multicasting algorithms, thus allowing for a new method of using multiple transmit antennas that delivers files of limited sizes with much reduced transmission delay.

30 [0026] The invented caching process is based on the idea that grouping defines, to a certain extent, the caching content, whereby the caching algorithm guarantees that some users store content that is highly related, while other users store content in such a way that there exists a certain degree of partial overlap.

35 [0027] Delivery of content is done by selecting some of the user groups at a time and creating a vector comprised of some of the desired content of each user in these selected groups. Specifically, by selecting some groups for content delivery, the transmitted vector becomes the combination (for example, the linear combination) of the precoded data vector for each group.

[0028] Decoding at the receiver side is done by using the cached content along with the channel state information of all receivers participating in the communication at that particular instance, meaning that only the receivers in the selected groups need to provide the channel state information.

40 [0029] In decentralized caching algorithms the delivery phase is done in a greedy manner, usually through graph coloring or some other methods. The here proposed method allows for an exponential decrease of the search space, because instead of having a graph where each graph node is a user, we have a graph where each graph node is a group.

Description of the drawings

45 [0030]

Figure 1 illustrates an architecture with a base station B (11), comprised of L antennas (12) and having access to a library of N files (13), serving K users (14), each with cache (15) of size equal to the size of M files.

50 Figure 2 illustrates one embodiment of a process leading to group-based caching at the receiving nodes.

Figure 3 illustrates one embodiment of a process of transmission of requested data. Delivery is done for K_g user groups at a time.

55 Figure 4 illustrates a process of decoding requested subfiles at receiving node.

Figure 5 shows a second embodiment corresponding to wired-network communications consisting of S transmitting servers (61), and K receiving users (66), via a network potentially consisting of routing nodes (63), collecting nodes

(65), and potentially combining nodes (64). Receivers request content from a library of N files. Each transmitter has a cache (62) of size equal to the size of M_T files, and each receiver has cache (67) of size equal to the size of M files.

Figure 6 illustrates an embodiment of a Caching process at transmitting servers in a wired network.

Figure 7 illustrates a process for group-based cache placement in wired network of the second embodiment.

Figure 8 illustrates a process of distributed transmission of requested data in wired network of the second embodiment. Delivery is done for K_g receiver-node groups at a time.

Figure 9 illustrates an embodiment of a process of decoding requested subfiles at a receiving node in the wired network of the second embodiment.

Figure 10 illustrates the D2D Model of a third embodiment.

Figure 11 illustrates the D2D Content Placement Process in the third embodiment.

Figure 12 illustrates the D2D Content Transmission Process in the third embodiment.

Figure 13 illustrates the Decoding Process in the D2D embodiment in the third embodiment.

Figure 14 illustrates the Distributed Based Stations of a fourth embodiment

Figure 15 illustrates the Content Placement at the Transmitters in the fourth embodiment.

Figure 16 illustrates the Content Placement at the Receivers in the fourth embodiment.

Figure 17 illustrates the Delivery of Content in the fourth embodiment.

Figure 18 illustrates the decoding Process at the receivers in the fourth embodiment.

Description of the preferred embodiments

[0031] There will now be described in details how one can improve content distribution in wireless and wired networks. The invention will be more particularly described in relation with a multi-antenna configuration and/or with multi-transmitter systems, either wireless or wired. The main purpose is to distribute digital content to receiving nodes which have storage memory and can cache part of this content and where there exists some interference as a result of the existence of many such receiving nodes with conflicting requested data.

[0032] While the current state-of-art methods exploit receiver side caches to achieve cache-aided multicasting (i.e., exploit the use of caches to allow for transmitting signals that are useful to many users at a time), such current state-of-art methods achieve multicasting gains that are severely limited, because of the well-known limitation that such gains require that each file be segmented into an exponentially large number of sub-files. Thus, given that the smallest segment can be that of a bit, the delivery of files with real-life size constraints incurs much limited multicasting gains and thus large delays. In the current state of art, this limitation exists (and is in fact made worse) when such state-of-art cache-aided multicasting methods are combined with multi-antenna precoding methods. Currently it is practically infeasible to meaningfully combine multi-antenna precoding methods with cache-aided multicasting methods, because of this aforementioned astronomical segmentation complexity brought about by existing state-of-art algorithms.

[0033] This limitation is tackled here by the design of new methods and apparatus which dramatically increases the state-of-art multicasting gains, by dramatically decreasing the aforementioned segmentation complexity and by dramatically decreasing the file size that is needed for high gains to appear.

[0034] In order to achieve such feasible solution, the inventors had to reconsider the formulation of the technical problem to solve. In summary, in the context of combining "MIMO" with "cache aided multicasting", it appears a huge problem which prohibits the design of any feasible solution. Indeed the key obstacle in the combination of MIMO and cache-aided multicasting is the entailed need to divide the files composing the library to be transmitted to the users into a myriad of small parts or packets. The technical problem to be solved is how to combine MIMO and cache-aided multicasting, in a way that allows for substantial multicasting gains, with a system that is still technically feasible.

[0035] By reconsidering such technical problem, the inventors have carried out a new approach, a new solution for achieving an exponential reduction in segmentation complexity of invented algorithm.

[0036] The current invention in this patent application here, manages - for the same exact problem of an L -antenna

base station serving K cache aided users - to reduce this complexity to a very manageable $\left(\frac{K}{L}\right)$ or even lower, closer

to $\left(\frac{1}{\gamma}\right)^{\frac{Ky}{L}-1}$. A good way to see the benefits of the current invention is with an example.

[0037] **Example:** Consider a base station that has $L = 20$ antennas, delivering files from a library of $N = 1000$ different movies to $K = 1000$ different wireless users, each having a storage cache of size being the equivalent to the size of $M = 20$ movies (i.e., each user can store 2% of the library). In this setting the best known multiplexing gain is $G_M = L = 20$

and the best known (theoretical) multicasting gain is $G_C = \frac{KM}{N} = 20$, thus the total gain (in theory) is $G_T = 20 + 20 = 40$. The state of art algorithms for combining multi-antenna precoding with cache-aided multicasting, in order to serve

the theoretical 40 users at a time, required segmentation complexity that was approximately $\binom{1000}{20} \times \binom{979}{19}$ which is approximately equal to 10^{68} (approaching the number of atoms in the universe). The invented algorithm in this patent application requires complexity approximately 100 (*one hundred*). Thus even if, with the state of art algorithms, we were allowed to segment a file into a million parts, this would allow for almost no cache-aided multicasting gains G_C , whereas the invented method here can achieve the theoretical multicasting gain $G_C = 20$ (adding an extra 20 users at a time per time-frequency slot, due to caching), with a very reasonable segmentation complexity of merely 100.

Technical solution: Key features that allow for the increased gains.

[0038] These gains are achieved mainly as a result of three innovative aspects that are key in our invention:

a) A novel caching (cache-placement) method which - unlike in any other cache-aided coded multicasting algorithms - regulates the amount of overlap between the contents of the caches, to be a variable that changes from cache to cache. This is a key ingredient in the success of the algorithm.

b) The introduction of the concept of user-grouping, which allows - in some instances of the problem - to set the caches of each group to be identical. This reduces the complexity of the problem without reducing the performance gains.

c) A novel transmission scheme that provides a total gain $G_T = L \times G'_C$ that is a multiplicative (rather than additive) combination of the multiplexing gain with the multicasting gain.

[0039] More particularly, there is provided a process that applies to a multi-antenna broadcast and multicast system, where a multi-antenna transmitting base station (B) transmits to receivers N_1, N_2, \dots, N_K which are equipped with caches. In this system, the process involves the steps of:

- grouping K users into groups containing a specified number of users per group,
- involving a first phase or time slot during which the base station employs two or more antennas to transmit information that is received and then partially cached at some or all of the caches of the receivers. The method of storing this information (i.e., the method of caching) is a function of at least the grouping of the users.
- involving a second phase or time slot during which some or all receiving nodes proceed to request additional information, e.g. the play of one particular movie, and based on this as well as on additional channel state information that the receivers feedback to the transmitter, the next function takes place:
- generating by the base station (B) of a multi-antenna transmit signal, typically in the form of a vector (or a matrix) whose length is equal to the number of transmit antennas, based on the users' grouping, with the purpose of transmitting the requested data files to said receiving nodes
- decoding by each participating node on the basis of the use of the information received in said first and second time slots and the information in the caches, so as to allow each receiving node N_k to compute an estimation of each requested file F_{rk}

[0040] A key aspect of the caching process above is based on the grouping of users which defines, to a certain extent, the caching content, whereby the caching algorithm guarantees that some users store content that is similar or identical,

while other users store content in such a way that there exists only a certain degree of partial overlap. The performance gains are in part because intra-group caches have very high redundancy (for example, caches of the same group can have all elements in common), while inter-group caches have reduced redundancy (in the sense that only a few files repeat in different caches). Having some users (members of the same group) with very similar or identical caches, manages to reduce the dimensionality of the multicasting problem, thus reducing complexity, while maintaining the overall gain.

[0041] This system combines multi-antenna precoding with cache-aided multicasting, to reduce the complexity associated with cache-aided multicasting algorithms¹, thus allowing for a new method of using multiple transmit antennas that delivers files (even of limited sizes), with much reduced transmission delay.

¹ Reducing complexity refers to, among other things, reducing the number of distinct transmissions which --- in the current state-of-art --- is astronomical.

[0042] In accordance with the invention described, the delivery of content is done by selecting some of the user groups at a time, and creating a vector comprised of some of the desired content of each user in these selected groups. Specifically, by selecting some groups for content delivery, the transmitted vector becomes the combination (for example, the linear combination) of the precoded data vector for each group.

[0043] For example, if in the absence of the invented algorithm, a group of users (group 1) could be served their desired data by a vector that is a result of multi-antenna precoding and then during another time slot another group of users (group 2) could be served similarly by another vector that results from a precoder, then the invented algorithm employs caching (cache-aided coded multicasting) techniques so that the two vectors are sent simultaneously (one "on top" of the other) thus serving both groups simultaneously.

[0044] Decoding at the receiver side is done by using the cached content along with the channel state information of all receivers participating in communication at that particular instance, meaning that only the receivers in the selected groups need to provide channel state information.

[0045] In decentralized caching algorithms the delivery phase is done in a greedy manner, usually through graph coloring or some other methods. The here proposed method allows for an exponential decrease of the search space, because instead of having a graph where each graph node is a user, we have a graph where each graph node is a group.

[0046] One thus sees that, by introducing here a new apparatus and method which precede the cache-placement and file delivery phases, we can improve the performance: from that of a system that achieves a compromised cache-aided multicasting gain $\overline{G}_C \ll G_C$, into the performance of a system that achieves in many cases of real-life interest, the entire (theoretical) multicasting gain $G_C \approx L \times \overline{G}_C$, thus boosting the previously compromised multicasting gains by a factor of up to L . This is the first time that multi-antenna methods are used to boost the real-life cache-aided multicasting gains. This boost is highlighted in the following example.

Example:

[0047] Consider a system with K users with caches, that achieves a real-life (severely compromised) multicasting gain of $\overline{G}_C = 6$. In a setting of a base station with $L = 5$ antennas, applying our method would boost this to an achieved multicasting gain of $L \times \overline{G}_C = 30$ (treating a total of $L + L \times \overline{G}_C = 35$ users at a time). Such gains have not been experienced before, to the best of our knowledge.

Note on the timeliness of the invention:

[0048] The above exploitation of multiple-antenna systems is consistent with today's trends, where more than one antenna in a single transmitter is becoming common even in domestic WiFi routers, while mobile vendors are moving toward large antenna arrays. Similarly, proposed solutions for 5G and beyond - which are focused on accommodating increased data traffic - aim to have tens or even hundreds of antennas. As a result, deploying the here proposed method is practical, and even with the use of a modest number of antennas, it can prove highly beneficial to decreasing the delays of content distribution. This invention also comes at a time when both, single-antenna caching and massive-antenna arrays without caching - when implemented independently - fail to scale so as to meet the current and future demands for content distribution, mostly due to their fundamental limitations. The herein proposed method combines these two in a complementary way which achieves to significantly slow-down their inevitable saturation in the regime of many users.

[0049] In order to illustrate the wide possibilities of application of the invention, there will now be described, in detail, a first embodiment covering a wireless configuration (I), a second embodiment directed to a wired network (II), a third embodiment describing a D2D (device to device) setting (III) and a fourth embodiment describing a wireless setting with distributed antennas (IV).

I. First embodiment (Wireless)

[0050] There will now be described a preferred first embodiment which shows the invented role given to the multi-antenna base station (B) and to the receiving nodes. The new approach combines multi-antenna precoding with cache-aided multicasting and manages to reduce the complexity associated with cache-aided multicasting algorithms and to allow for much reduced transmission delay when delivering files of any given sizes.

[0051] The preferred embodiment, Figure 1, corresponding to a base station B 11, comprised of L antennas 12, and having access to a library of N files 13, serving K users 14, each fitted with a cache 15 storage of size equal to the size of M files.

[0052] In this first embodiment, there are three basic processes embodying the two time slots: the process, of the first time slot, that leads to content placement in the caches of the receiving nodes (Figure 2), the process, of the second time slot, that leads to transmission of content by the base station B (Figure 3) and the process of decoding at the receiving nodes (Figure 4).

Process that leads to content placement in the caches of the receiving nodes:

[0053] As illustrated in figure 2, the first process comprises of the following steps 21-26 in sequence to be performed during the first time-slot.

[0054] Steps 21 and 22 are respectively "user enumeration" and "evaluation of grouping function and user grouping".

With respect to these steps, users are enumerated and assigned into one of the $\frac{K}{L}$ groups. The grouping function will take as input, among other things, the user enumeration, the network topology (link strengths), as well as the cache capabilities of the different receiving nodes. In this specific embodiment here, where the link strengths are the same for each user, and where the cache sizes are also the same across the different receivers, the grouping function is simple

and it assigns each group with L users, as follows: users $\left\{1, \frac{K}{L} + 1, \frac{2K}{L} + 1, \dots, \frac{(L-1)K}{L} + 1\right\}$ are placed in group

1, users $\left\{2, \frac{K}{L} + 2, \frac{2K}{L} + 2, \dots, \frac{(L-1)K}{L} + 2\right\}$ are placed in group 2, and so on.

[0055] The next step 23 is "Creation of subpacketization indices", which results in the calculation of subpacketization indices τ . To do this, we gather a set T that is some subset of S^* sets τ , where each set τ consists of Q distinct numbers

from $\left\{1, 2, \dots, \frac{K}{L}\right\}$. In some specific embodiments, Q can take the value $Q = \frac{K\gamma}{L}$.

[0056] A next step 24, so-called "File segmentation into subpackets", consists of file segmentation, wherein each file $F_n (n = 1, 2, \dots, N)$ is segmented into S^* subfiles $F_{n,\tau}$ where each set τ in T is used to label a different segment of the file.

[0057] In the next - so called "Transmission of file library" - step 25, the base station (B) sequentially transmits x_1 which maps information from F_1, F_2, \dots, F_N , and each receiving node $N_k (k = 1, 2, \dots, K)$ processes their received signal $f_k(x_1)$ to get an estimate of different parts of the library F_1, F_2, \dots, F_N . The transmission is multicast or partial unicast and it can employ one or more transmit antennas.

[0058] At the completion of step 25, the process described in figure 2 then proceeds to a step 26 - designated "Group-based content cache placement" - which consists in group-based caching, where each user, depending at least on its group, caches so that the cache overlap - corresponding to the amount of common subfiles - between any two caches, is calculated to be at least a function of whether the two caches belong in the same group or not. In this particular embodiment where the cache sizes are of the same size, and where the link strengths are statistically similar, the caching

is as follows: Each user belonging to group i , (where $i = 1, 2, \dots, \frac{K}{L}$), caches every subfile $F_{n,\tau}$ as long as i can be found inside the set τ .

Process that leads to transmission of content by the base station B:

[0059] With respect to figure 3, there will now be described the second process to be performed in the second time slot, consisting of a sequence of steps 31-39, for gathering file requests (step 31), then creating the transmission indices for active groups (step 32), then identifying K_g active groups (step 33), then gathering channel estimates for K_g active groups (step 34), then creating the precoding vector for an active group (step 35), repeating for each active group (step

36), then combining the precoding vectors of the K_g active groups (step 37), and then transmitting the composite vector simultaneously to all K_g active groups (step 38). The process (32)-(38) is repeated for different sets of K_g active groups (step 39).

[0060] More particularly, the process of figure 3 starts with a first step 31 being the gathering of file requests, where the receiving nodes request from the transmitter to deliver specific files from the library. In this first preferred embodiment, each receiving node N_k requests one file, by communicating to the base station B the information that describes the label of their requested file F_{r_k} ($r_k = 1, 2, \dots, N$).

[0061] Then, the process proceeds to a step 32, consisting of gathering transmission indices. In this step, indices χ are created that will describe which $Q + 1$ groups will be active, which in turn describes which $Q + 1$ groups to transmit to at any given time. To do this, the process gathers a set X that is some subset of possible sets χ , where each set χ

consists of $Q + 1$ distinct numbers from $\left\{1, 2, \dots, \frac{K}{L}\right\}$.

[0062] Then, after completion of step 32, the process proceeds to a step 33, which is about identifying the active group. Here, sets χ are sequentially picked from X . Once a set χ is picked, this immediately identifies the $Q + 1$ active groups, which include group i as long as i can be found in the set χ . Thus each χ automatically identifies the $L \times (Q + 1)$ so

called 'active receiving nodes'. In some specific embodiments, $Q + 1$ can take the value $Q + 1 = \frac{K\gamma}{L} + 1$. Given a χ , the step identifies $(Q + 1) \times L$ active receiving nodes.

[0063] Then, in a step 34, the process performs the gathering of channel estimates for the K_g active groups, where each active receiving node N_k communicates their channel state information h_k to the transmitter. As known by a skilled man, this feedback can also take other forms, that include receiver state, precoder preference, location information, etc.

[0064] In a step 35, the process creates a precoding vector for an active group, by precoding L subfiles for L users in the active group. Here, for a given active group i , the transmitter generates the L -length signal vector

$$\mathbf{w}_{G_i}(\chi) \cdot (H_{G_i})^\perp$$

where $(H_{G_i})^\perp$ is the precoding matrix relative to the channel of group i , and where

$$\mathbf{w}_{G_i}(\chi) = [F_{r_{G_i(1),\chi^i}}, F_{r_{G_i(2),\chi^i}}, \dots, F_{r_{G_i(L),\chi^i}}]$$

is the vector comprised of the vector representation of the subfiles $F_{r_{G_i(1),\chi^i}}, F_{r_{G_i(2),\chi^i}}, \dots, F_{r_{G_i(L),\chi^i}}$ wanted by the L users of group i .

[0065] Then, in a step 36, the process repeats the above step (35) for each active group.

[0066] In a step 37, the process combines the precoding vectors $\mathbf{w}_{G_i}(\chi) \cdot (H_{G_i})^\perp$ of the K_g active groups, to form a composite vector

$$\mathbf{x}_{t_s} = \sum_{i \in \chi} \mathbf{w}_{G_i}(\chi) \cdot (H_{G_i})^\perp$$

where the summation is over integer numbers i that comprise the currently chosen set χ of active groups.

[0067] In a step 38, the base-station B transmits the composite vector \mathbf{x}_{t_s} of length L , simultaneously to all K_g active groups. Each scalar of the vector corresponds to one antenna of the transmitter.

[0068] At last, in a step 39, the process repeats the steps 32 - 38 for different sets of K_g active groups.

Process of decoding at the receiving nodes:

[0069] With respect to figure 4, there will now be described the third process - which is presented here from the point of view of any one user - for decoding at the receiving nodes and which comprises the sequence of steps 41-49 in

sequence.

[0070] The third process starts with a step 41, designated "Is user in set of active groups?", wherein a test is performed to determine whether the user belongs to the upcoming transmission, which means that the receiver is active.

[0071] The result of the test is processed in a step 42.

5 [0072] If the user does not belong to the upcoming transmission, the process proceeds to a step 48 - designated "Repeat for new set of K_g until file is received" - ,wherein the receiver will wait until the next transmission.

[0073] If the user belongs to the upcoming transmission, the process proceeds to a step 43, designated "Receiver broadcast its channel information", and the receiver will broadcast its channel state information.

10 [0074] Then, in a step 44 designated "Receiver records channel information from other active users", the receiver will listen for and record the other channel information from the other active users.

[0075] Then, the process proceeds to a step 45 designated "User receives signal meant for K_g active groups", wherein the receiver records the received signal which is meant for the active groups.

15 [0076] In a step 46 - designated "Cache-out" intergroup interfering messages based on group-based multicasting scheme" - the receiver uses channel information and cached content to remove (cache-out) inter-group interference. This is achieved by using a state-of-art cache-aided multicasting method well known to a skilled man, and which is

elevated in this step to a group-level, such that, by design, the sets T and X , guarantee that all the elements of $w_{G_i}(X) =$

20 $[F_{r_{G_i(1),X \setminus i}}, F_{r_{G_i(2),X \setminus i}}, \dots, F_{r_{G_i(L),X \setminus i}}]$ are cached at the rest participating (active) groups' caches, except the one that desires them (i.e., except at the caches of group i). This means that each receiver of a specific active group holds in its cache all the file segments that have been requested by the users of the other groups that are active at that particular instance.

25 [0077] In a step 47 - designated "Null-out" intra-group interference, Decode own subfile using precoding-based decoder", the receiver applies a standard decoder that reflects the employed state-of-art precoder, to receive its desired subfile. By design of the precoder-and-decoder (which for a specific embodiment can be a simple ZF precoder), a message decoded by a receiver does not contain any file segment intended for another receiver of the same group.

[0078] By combining steps (46) and (47), the receiver can remove all other Q unwanted subfiles, and decode its own subfile.

30 [0079] Step (48) consists of repeating the process (43) to (47), until the entire desired file is received.

EXAMPLE BASED ON FIRST EMBODIMENT

35 [0080] Assume a base station B with $L = 4$ transmit antennas is deployed in a mall. Let there be a library of $N = 20$ files (for example $N = 20$ different movies) and let there be $K = 20$ wireless users (previously referred to as receiving nodes) labeled by $1, 2, \dots, 20$. Let each receiving node $k = 1, 2, \dots, 20$ have a cache (storage) of size corresponding to M

$= 8$ (equivalent to the size of 8 files), thus corresponding to a normalized cache size of $\gamma = \frac{M}{N} = \frac{8}{20} = \frac{2}{5}$.

40

1. Phase 1: Placement (first time slot)

45

- *User enumeration and grouping:* Users are divided into 5 groups of size $L = 4$ users. The groups are $G_1 = \{1, 6, 11, 16\}$ (Group 1), $G_2 = \{2, 7, 12, 17\}$ (Group 2), $G_3 = \{3, 8, 13, 18\}$ (Group 3), $G_4 = \{4, 9, 14, 19\}$ (Group 4) and $G_5 = \{5, 10, 15, 20\}$ (Group 5).

$$S^* = \binom{\frac{K}{L}}{\frac{Ky}{L}} = \binom{5}{2} = 10 \text{ different sets } \tau,$$

50

- *Subpacketization indices:* We create the set T that includes $Q = \frac{Ky}{L} = 2$ distinct numbers from $\{1, 2, \dots, 5\}$. Thus $T = \{(1, 2), (1, 3), (1, 4), (1, 5), (2, 3), (2, 4), (2, 5), (3, 4), (3, 5), (4, 5)\}$.

55

- *File segmentation:* Each file $F_n (n = 1, 2, \dots, N)$ is segmented into $S^* = 10$ subfiles $F_{n,\tau}$ where each set τ in the above T is used. For example, file F_1 is segmented into subfiles $\{F_{1,(1,2)}, F_{1,(1,3)}, F_{1,(1,4)}, \dots, F_{1,(4,5)}\}$ and similarly with the rest of the files F_2, \dots, F_{20} .

- *First phase transmission:* The base station (B) sequentially transmits x_1 which maps information from F_1, F_2, \dots, F_N .

- *Group-based caching:* Based on x_1 , cache-placement commences, so that each user belonging to Group i

$$\left(i = 1, 2, \dots, \frac{K}{L}\right)$$

caches every subfile $F_{n,\tau}$ as long as i can be found inside the set τ . The cache placement is defined by the following: $Z_1 = \{(1,2), (1,3), (1,4), (1,5)\}$ (cache for Group 1), $Z_2 = \{(1,2), (2,3), (2,4), (2,5)\}$ (cache for Group 2), $Z_3 = \{(1,3), (2,3), (3,4), (3,5)\}$ (cache for Group 3), $Z_4 = \{(1,4), (2,4), (3,4), (4,5)\}$ (cache for Group 4) and $Z_5 = \{(1,5), (2,5), (3,5), (4,5)\}$ (cache for Group 5). This means that, for example, all users of group 1 (i.e., users 1,6,11,16) will cache all subfiles $F_{n,(1,2)}, F_{n,(1,3)}, F_{n,(1,4)}, F_{n,(1,5)}$, for each file F_n . Similarly all users of group 2 (i.e., users 2,7,12,17) will cache all subfiles $F_{n,(1,2)}, F_{n,(2,3)}, F_{n,(2,4)}, F_{n,(2,5)}$, again for each file F_n , and so on. It can be seen that inter-group caches have 1 out of 5 parts in common (e.g. Group 1 and Group 2 both have segments labeled by (1,2)), while intra-group caches have all elements in common (i.e., caches of the same group are identical).

2. Phase 2: Delivery (second time slot)

- *File requests:* Each user asks for any one file from the library, so user 1 asks for file F_{r_1} , user 2 for F_{r_2} , and so on, where $r_k = 1, 2, \dots, 20$.

- *Transmission indices and active group identification:* We consider the set $X =$

$\{(1,2,3), (1,2,4), (1,2,5), (1,3,4), (1,3,5), (1,4,5), (2,3,4), (2,3,5), (2,4,5), (3,4,5)\}$ consisting of sets χ , where each set χ consists of $Q + 1 = 3$ distinct numbers from the set $\{1,2,3,4,5\}$. Each different χ describes a different set of $Q + 1 = 3$ active groups, corresponding to $(Q + 1) \times L = 12$ active users. For example, when $\chi = (1,2,3)$, it means that the active groups are groups 1,2 and 3, which are the groups that the base station will transmit to during this time.

- *Feedback to transmitter and Transmission to groups:* The first transmission is denoted by x_{123} , where the index implies that all users belonging in groups 1,2 and 3 are going to be receiving a desired subfile. Each active receiving node N_k (i.e., any user from groups 1,2 or 3) communicates channel state information h_k back to the transmitter. The transmitted vector, from the base station B to the active receivers then takes the form

$$x_{123} = (H_{G_1})^{-1} \cdot \begin{pmatrix} F_{r_{1,(2,3)}} \\ F_{r_{6,(2,3)}} \\ F_{r_{11,(2,3)}} \\ F_{r_{16,(2,3)}} \end{pmatrix} + (H_{G_2})^{-1} \cdot \begin{pmatrix} F_{r_{2,(1,3)}} \\ F_{r_{7,(1,3)}} \\ F_{r_{12,(1,3)}} \\ F_{r_{17,(1,3)}} \end{pmatrix} + (H_{G_3})^{-1} \cdot \begin{pmatrix} F_{r_{3,(1,2)}} \\ F_{r_{8,(1,2)}} \\ F_{r_{13,(1,2)}} \\ F_{r_{18,(1,2)}} \end{pmatrix}$$

[0081] Then, the received messages - focusing on the case of group 1 ($G_1 = \{1,6,11,16\}$) - stacked as a vector to incorporate the received signals at users $\{1,6,11,16\}$, takes the form:

$$y_{123}(G_1) = \begin{pmatrix} F_{r_{1,(2,3)}} \\ F_{r_{6,(2,3)}} \\ F_{r_{11,(2,3)}} \\ F_{r_{16,(2,3)}} \end{pmatrix} + H_{G_1} \cdot (H_{G_2})^{-1} \cdot \begin{pmatrix} F_{r_{2,(1,3)}} \\ F_{r_{7,(1,3)}} \\ F_{r_{12,(1,3)}} \\ F_{r_{17,(1,3)}} \end{pmatrix} + H_{G_1} \cdot (H_{G_3})^{-1} \cdot \begin{pmatrix} F_{r_{3,(1,2)}} \\ F_{r_{8,(1,2)}} \\ F_{r_{13,(1,2)}} \\ F_{r_{18,(1,2)}} \end{pmatrix}$$

[0082] As seen from the above equation, each user in Group 1 receives the desired subfile plus a linear combination of all the subfiles intended for the users in the other participating groups (groups 2 and 3 in this case), where this linear combination depends on the channels of these $K\gamma + L$ participating (active) users. By using the cached content and the knowledge of the channel coefficients, each user can decode its desired subfile.

[0083] In the same manner, we proceed with the rest of the transmissions, i.e., with the rest of the χ in X , i.e., with $\chi = (1,2,4)$ corresponding to active groups 1,2 and 4, then with $\chi = (1,2,5)$ corresponding to active groups 1,2 and 5, and so on.

[0084] From the above algorithm we can see that there are $K\gamma + L$ users served at a time, while the number of file

$$S^* = \binom{\frac{K}{L}}{\frac{K\gamma}{L}} = \binom{5}{2} = 10.$$

segments (subpacketization) is needed to achieve this same performance would have been

$$S^* = \binom{K}{K\gamma} \times \binom{K-K\gamma-1}{L-1} = \binom{20}{8} \times \binom{11}{3} = 125970 \times 165 = 20878050.$$

The complexity reduces from approximately twenty million (distinct transmissions), to just ten distinct transmissions.

[0085] Another interpretation of the above is that for any state-of-art placement-delivery which is constrained by the file size, the here proposed method can achieve a boost of L times on (cache-aided multicasting gains). Specifically, if a cache-aided multicasting algorithm can achieve multicasting gain of \overline{G}_C for some maximum allowable number F of segments per file (maximum subpacketization level of F) then the here proposed algorithm can achieve a multicasting gain of $\min\{L \times \overline{G}_C, K\gamma + 1\}$ under the same segmentation (complexity) constraint F .

II. Second embodiment (Wired network)

[0086] There will now be described a second embodiment which shows the invented role given to multiple servers and to the receiving nodes, in a wired network. The new approach combines the previously described technique of fusing multi-transmitter precoding with cache-aided multicasting, to again reduce the complexity of cache-aided multicasting algorithms and to allow reduced transmission delays when delivering files of any given sizes, from a set of transmitting servers to a set of receiving nodes.

[0087] Figure 5 more particularly illustrates such second embodiment, showing S transmitting servers **61**. As before, there is a library of N files, but each transmitting server only has access to the part of the library that is stored at the cache of that transmitter **62**. Each transmitter's cache **62** has size equal to the size of M_T files. The transmitting servers aim to deliver data to K receiving users **66**, each with cache **67** of size equal to the size of M files. The medium of communication is a wired network, potentially consisting of routing nodes **63** and collecting nodes **65**, as well as potentially consisting of combining nodes **64** which can combine receiving inputs at the node, to produce a composite output from that node.

[0088] There are four distinct processes: the process that leads to content placement in the caches of the receiving nodes (Figure 6), the process that leads to content placement in the caches of the receiving nodes (Figure 7), the process that leads to transmission of content by the servers (Figure 8) and the process of decoding at the receiving nodes (Figure 9). One essential difference in this wired embodiment is that, unlike in the wireless embodiment, which considered signals and channel-transfer functions to be consisting of real or complex numbers, here the network-transfer functions and the signals consist of elements from a finite field. By design, the size of the finite fields can be chosen to be large enough, so that the network transfer matrices can be of higher rank.

[0089] Process that leads to content placement in the caches of the transmitting servers: The first process, Figure 6, starts with a step **71** consisting of server-enumeration, where the servers are labeled as $s = 1, 2, \dots, S$. The second step **72**, designated "Pick server $s: s = 1, 2, \dots, S$ ", initiates the placement algorithm for the transmitter-side caches.

For fractional cache of $\gamma_T = \frac{M_T}{N}$, the algorithm places, in a step **73**, any subfile at exactly $S \times \gamma_T$ transmitting servers. These $S \times \gamma_T$ transmitting servers - for that specific subfile - can play the previous role of the $L = S \times \gamma_T$ transmit antennas of the single base station B. The algorithm consecutively caches entire files F_1, F_2, \dots, F_N into the transmitters, such that the first transmitter caches the first M files, the second transmitter the next M files, and so on, modulo N . Using $Z_{Tx,s}$ to denote the cache of transmitting server $= 1, 2, \dots, S$, then

$$Z_{Tx,s} = \{ F_{1+(n-1) \bmod N} : n \in \{1 + (s-1)M, \dots, Ms\} \}.$$

[0090] This guarantees the memory constraints and the aforementioned requirement that each subfile resides in $S \times \gamma_T$ transmitters.

[0091] Later, in the transmission of any given subfile, these $S \times \gamma_T$ transmitters, with knowledge of the network topology (i.e., with network of the equivalent of CSIT, corresponding to the network transfer matrix), can now play the role of the aforementioned $L = S \times \gamma_T$ antennas, and thus can precode this said subfile using the exact same precoders described before, allowing for simultaneous separation of the $L = S \times \gamma_T$ streams within any given group G_g of $L = S \times \gamma_T$ receivers.

As before, the aforementioned caching allows for treatment of $\frac{K}{L}\gamma + 1$ groups at a time and a treatment of $S \times \gamma_T + K \times \gamma \leq K$ users at a time.

[0092] The process is terminated, in a step **74**, when all caches are full ($s = S$), and when each file appears in $L = S \times \gamma_T$ servers.

[0093] Process that leads to content placement in the caches of the receiving nodes: This process, illustrated in figure 7, consists of the process of user-enumeration (step **81**), calculation of grouping function and user grouping (step **82**), calculation of the subpacketization indices (step **83**), file segmentation (step **84**), sequential transmission of content from the library of N files (step **85**) and group-based cache placement (step **86**). In the following we will set $L = S \times \gamma_T$ which will make some of the steps here similar to the steps of the first preferred embodiment where, in a wireless network, a base station B with L antennas was serving K receiving nodes.

[0094] With respect to the step of user enumeration (step **81**) and user grouping (step **82**), users are enumerated and assigned into groups as before. The grouping function will take as input, among other things, the user enumeration, the network topology and network connectivity, the role assigned to the combining nodes (**64**), as well as the cache capabilities of the different receiving nodes. Taking these parameters into consideration, the grouping function will be part of a basic optimization for decreasing the overall delay. In this specific embodiment here - where the link strengths are statistically the same for each user (the links have similar capacity in the long term) and where the cache sizes are also the same across the different receivers - the grouping function assigns each group with L users, as follows: users

$\left\{ 1, \frac{K}{L} + 1, \frac{2K}{L} + 1, \dots, \frac{(L-1)K}{L} + 1 \right\}$ in group 1, users $\left\{ 2, \frac{K}{L} + 2, \frac{2K}{L} + 2, \dots, \frac{(L-1)K}{L} + 2 \right\}$ in group 2, and so on.

[0095] The next step **83** is to calculate the subpacketization indices τ . This is done as before: a set T is created that

is some subset of S^* sets τ , where each set τ consists of Q distinct numbers from $\left\{ 1, 2, \dots, \frac{K}{L} \right\}$. In some specific

embodiments, Q can take the value $Q = \frac{K\gamma}{L}$.

[0096] The next step **84** is file segmentation, where each file $F_n (n = 1, 2, \dots, N)$ is segmented into S^* subfiles $F_{n,\tau}$ where each set τ in T is used to label a different segment of the file.

[0097] In the next step **85**, the S servers sequentially transmit part of their content: each server $s = 1, 2, \dots, S$ transmits - one after the other - signals x_s which map information from F_1, F_1, \dots, F_N . At the same time, each receiving node $N_k (k = 1, 2, \dots, K)$ then processes their received signal $f_k(x_1, x_2, \dots, x_S)$ to get an estimate of different parts of the library F_1, F_2, \dots, F_N . The transmission x_s at each server can account for any previously transmitted content by the servers and it can account for network connectivity.

[0098] The next step **86** is the step of group-based caching, where each user, depending at least on its group, caches so that the cache overlap - corresponding to the amount of common subfiles - between any two caches, is calculated to be at least a function of whether the two caches belong in the same group or not. In this particular embodiment where the cache sizes are of the same size, and where the link strengths are statistically similar, the caching is as follows:

Each user belonging to group i , (where $i = 1, 2, \dots, \frac{K}{L}$), caches every subfile $F_{n,\tau}$ as long as i can be found inside the set τ .

[0099] Process that leads to transmission of content by the S transmitting servers: The third process, Figure 8, consists of steps **91-100** for gathering file requests (step **91**), then creating the indices for active groups of receiving users (step **92**), then identifying K_g active groups of receiving users (step **93**), then identifying $L = S \times \gamma_T$ active transmitting servers (step **94**), then gathering network connectivity information for K_g active groups (step **95**), then creating the precoding vector for an active group (step **96**), repeating for each of the K_g currently active groups (step **97**), then combining the precoding vectors of the K_g currently active groups (step **98**) and then transmitting the composite vector simultaneously to all K_g currently active groups (step **99**). The process (92)-(99) is repeated for different sets of K_g active groups (step **100**).

[0100] Step **91** starts the process with the action of file requests, where the receiving nodes request delivery of specific files from the library. In this preferred embodiment, each receiving node N_k requests one file, by communicating to the transmitting servers the information that describes the label of that receiving node's requested file $F_{r_k} (r_k = 1, 2, \dots, N)$. Furthermore, in this preferred embodiment, the file requests are treated concurrently.

[0101] Step **92** consists of creating indices for active receiver-node groups. In this step, indices χ are created that will describe which $Q + 1$ receiver-node, groups will be simultaneously active, which in turn describes which $Q + 1$ groups

to transmit at any given time. To do this, we gather a set X that is some subset of possible sets χ , where each set χ

$$\left\{1, 2, \dots, \frac{K}{L}\right\}.$$

consists of $Q + 1$ distinct numbers from

5 **[0102]** Step 93 is about identifying the active receiver groups, among the subset of all possible groups. Here, sets χ are sequentially picked from X . Once a set χ is picked, this immediately identifies the $Q + 1$ active groups, which include group i as long as $i \in \chi$, i.e., as long as the number i can be found in the set χ . Thus, each χ automatically identifies the $L \times (Q + 1)$ so called "active receiving nodes". In some specific embodiments $Q + 1$ can take the value

$$10 \quad Q + 1 = \frac{KY}{L} + 1.$$

Given a set χ , the step identifies $(Q + 1) \times L$ active receiving nodes.

[0103] Step 94 consists in the identification of $L = S \times \gamma_T$ active transmitting servers. Recalling from the process (step 72) of Figure 6, the transmitter-side caching step guarantees the requirement that each subfile resides in $L = S \times \gamma_T$ transmitters. Consequently, given the set of requested files from step 91, and given the set of $(Q + 1) \times L$ active receiving nodes from step 93, this step here 94, applies the transmitter-side placement formula $Z_{T_x, s} = \{F_{1+(n-1) \bmod N}, n \in \{1 + (s - 1)M, \dots, Ms\}\}$ to identify the $L = S \times \gamma_T$ transmitters that will deliver the requested subfiles.

15 **[0104]** Step 95 is about gathering the network connectivity and network transfer function information corresponding to the network between the K_g active receiver-node groups (from step (93)) and the associated L transmitting servers (from step (94)), such feedback can also take forms that include receiver state, precoder preference, location information, etc. Using basic techniques, this information - which describes the state of the wired communication medium, and which accounts for connectivity, attenuation coefficients, as well as for the actions of the combining nodes - is communicated among the receiving nodes, and to the currently transmitting servers.

20 **[0105]** Step 96 creates a distributed precoding vector for an active receiver-node group, by precoding L subfiles, across the L different transmitters, for L users in the active group. Here, for a given active group i , the associated L servers jointly generate (using knowledge of each other's state information) the L -length signal vector

$$w_{G_i}(\chi) \cdot (H_{G_i})^\perp$$

30 where $(H_{G_i})^\perp$ is the precoding matrix relative to the channel between group i and the associated L servers and where

$$35 \quad w_{G_i}(\chi) = [F_{r_{G_i(1), \chi \setminus i}}, F_{r_{G_i(2), \chi \setminus i}}, \dots, F_{r_{G_i(L), \chi \setminus i}}]$$

is the vector comprised of the vector representation of the subfiles $F_{r_{G_i(1), \chi \setminus i}}, F_{r_{G_i(2), \chi \setminus i}}, \dots, F_{r_{G_i(L), \chi \setminus i}}$ wanted by the L users of group i . As stated before, one essential difference here is that $w_{G_i}(\chi)$, H_{G_i} and by extension $(H_{G_i})^\perp$ and $w_{G_i}(\chi) \cdot (H_{G_i})^\perp$ consist of elements from a finite field, unlike in the wireless embodiment which considers channels and signals consisting of real or complex numbers. By design, the size of the finite fields can be chosen to be large enough, so that the transfer matrices H_{G_i} can be of higher rank.

40 **[0106]** Step 97 repeats the above step (96) for each active group.

45 **[0107]** Step 98 combines the precoding vectors $w_{G_i}(\chi) \cdot (H_{G_i})^\perp$ of the K_g active groups, to form a composite vector

$$50 \quad x_{t_s} = \sum_{i \in \chi} w_{G_i}(\chi) \cdot (H_{G_i})^\perp$$

where the finite-field summation is over integer numbers i that comprise the currently chosen set χ of active groups.

[0108] In step 99 the L currently chosen servers jointly transmit the composite vector x_{t_s} of length L , simultaneously to all K_g active groups. Each scalar of the vector corresponds to one server.

55 **[0109]** Step 100 consists of repeating the process (92)-(99) for different sets of K_g active receiver-node groups and potentially different sets of active servers.

[0110] Process of decoding at the receiving nodes: In this embodiment, the steps of the fourth process for decoding

are presented in Figure 9, steps **101-109**. The process - which is presented here from the point of view of any one receiving user - consists of the step where the user is notified or calculates (step **101**) if it belongs in the upcoming transmission, which means that the receiver is active. If NO (step **102**), the receiver will wait until the next transmission to recheck (step **106**). If YES (step **102**), the receiver will broadcast it's network state information (step **103**), and then it will listen for and record (step **104**) the other channel information from the other active users. Then the receiver records the received signal (step **105**), which is meant for the active groups.

[0111] In step (**106**) the receiver uses network-state information and cached content to remove (cache-out) inter-group interference. This is achieved by using a state-of-art cache-aided multicasting method, which is elevated in this step to

a group-level such that it guarantees that all the elements of $\mathbf{w}_{G_i}(\chi) = [F_{r_{G_i(1),\chi^i}}, F_{r_{G_i(2),\chi^i}}, \dots, F_{r_{G_i(L),\chi^i}}]$ are cached at the rest participating (active) groups' caches, except the one that desires them (i.e., except at the caches of group i). This means that each receiver of a specific active group holds in its cache all the file segments that have been requested by the users of the other groups that are active at that particular instance.

[0112] In step **107** the receiver applies a standard decoder that reflects the employed state-of-art precoder, to receive its desired subfile. By design of the precoder-and-decoder, which for a specific embodiment can be a simple ZF precoder, a message decoded by a receiver does not contain any file segment intended for another receiver of the same group.

[0113] By combining steps **46** and **47**, the receiver can remove all other Q unwanted subfiles and decode its own subfile.

[0114] Step **108** consists of repeating the process (steps 43 to 47), until the entire desired file is received.

III. Third embodiment (D2D)

[0115] There will now be described a preferred third embodiment which shows the invented role given to wireless nodes sharing the same medium, equipped with a cache and wanting to exchange data from some library of pre-defined content. The new approach leverages cooperation between the nodes to overcome subpacketization constraints. In theory, in a subpacketization unconstrained system one can multicast a message to $K\gamma$ other nodes at a time, where K

is the number of nodes and $\gamma = \frac{M}{N}$ is the fractional capacity of the cache of a user, meaning that each user is storing the equivalent of M files out of a library of N files. In practice, though, these gains are severely limited by the aforementioned required subpacketization.

The new approach invented here manages to satisfy the demands of $K\gamma$ users at a time, by employing a first-of-its-kind node cooperation that dramatically reduces the subpacketization requirements and, thus, dramatically boosts the rate-performance of the system. The invention involves the cooperation of $L \leq K\gamma$ nodes at a time, were these cooperating nodes create a distributed precoded vector with messages for $K\gamma$ other nodes, which are served simultaneously with reduced subpacketization. The invented approach allows for a small number of cooperating nodes, but offers exponential reductions in the subpacketization costs; while still achieving the full gain of serving $K\gamma$ users at a time.

[0116] The preferred embodiment, illustrated in Figure 10, corresponds to K independent nodes N_1, N_2, \dots, N_K (**131**), each with a single antenna (**132**), and each equipped with a cache (**133**) of size equal to the size of M files while each user can ask for any one file from a library of N files.

[0117] In this third embodiment, there are three basic processes embodying two time slots: i) the process, of the first time slot, that leads to content placement in the caches of the nodes (Figure 11), ii) the process of the second time slot that leads to transmission of content by the nodes (Figure 12), and iii) the process of decoding at the nodes (Figure 13).

Process that leads to content placement in the caches of the nodes:

[0118] As illustrated in Figure 11, the first process comprises the following steps **111-114** in sequence to be performed during the first time-slot. In step (**111**) the number of cooperating nodes L is to be decided. One such L could be the

smallest integer satisfying the subpacketization constraint $\frac{K\gamma}{L} \left(\frac{K/L}{K\gamma/L} \right) \leq F^*$, where F^* is the maximum allowed sub-

packetization. Then in step (**112**), users are assigned into one of the $\frac{K}{L}$ groups, for example as follows: users

$\left\{ 1, \frac{K}{L} + 1, \frac{2K}{L} + 1, \dots, \frac{(L-1)K}{L} + 1 \right\}$ are placed in group 1, users $\left\{ 2, \frac{K}{L} + 2, \frac{2K}{L} + 2, \dots, \frac{(L-1)K}{L} + 2 \right\}$ are

placed in group 2, and so on.

[0119] The next step **113** is "File Segmentation" which starts with the calculation of subpacketization indices p and z .

To do this, we gather the set $P = \left\{1, 2, \dots, \frac{K\gamma}{L}\right\}$ and the set T that is a set comprised of subsets, denoted by τ , of Q

$$Q = \frac{K\gamma}{L}.$$

distinct elements from the set $\left\{1, 2, \dots, \frac{K}{L}\right\}$. In some embodiments, Q can take the value

[0120] In the file segmentation, wherein each file F_n ($n \in \{1, 2, \dots, N\}$) is segmented into S^* subfiles $F_{n,p,\tau}$ where each subset $(p, \tau) \in P \times T$ is used to label a different segment of the file.

[0121] At step 114, named "Group-based content cache placement", where content is saved in each user according

to the group a user belong to. For example, users in Group i store all files $F_{n,p,\tau} \forall n \in [N], \forall p \in \left[\frac{K\gamma}{L}\right], i \in \tau$.

Process that leads to transmission of content by the nodes:

[0122] With respect to Figure 12, there will now be described the second process to be performed in the second time slot, consisting of a sequence of steps 121-129, for gathering file requests (step 121), then, each group in sequence

(122) selects some groups, for example $\frac{K\gamma}{L}$ groups (123) gets channel state information from these groups (124), creates a distributed precoder for each group (125), precodes messages for each group using these precoders (126), combines the precoded messages (127) and repeats the process for a new set of groups (128). Finally, the process is repeated using a next group as a transmitter.

Process of decoding at the receiving nodes:

[0123] With respect to figure 13, there will now be described the third process - which is presented here from the point of view of any one user - for decoding at the receiving nodes and which comprises the sequence of steps 141-149 in sequence.

[0124] The third process starts with a step 141, designated "Is user in set of active groups?", wherein a test is performed to determine whether the user belongs to the upcoming transmission, which means that the receiver is active.

[0125] The result of the test is processed in a step 142.

[0126] If the user does not belong to the upcoming transmission, the process proceeds to step 148 - designated "Repeat for new set of K_g until file is received", wherein the receiver will wait until the next transmission.

[0127] If the user belongs to the upcoming transmission, the process proceeds to a step 143 designated "Recover broadcasts its channel information", where the receiver will broadcast its channel state information.

[0128] Then, in step 144 designated "Receiver records channel information from other active users", the receiver will listen for and record channel information from the other active users.

[0129] Then, the process proceeds to a step 145 designated "User receives signal meant for K_g active groups", wherein the receiver records the received signal, which is meant for the active groups.

[0130] In a step 146 designated "Cache-out" intergroup interfering messages based on group-based multicasting scheme", the receiver uses channel information and cached content to remove (cache-out) inter-group interference. This is achieved by using a state-of-art cache-aided multicasting method well known to a skilled man, and which is elevated in this step to a group-level such that, by design, the sets T and X , guarantee that all the elements of

$w_{G_i}(X) = [F_{r_{G_i(1),j,X|i}}, F_{r_{G_i(2),j,X|i}}, \dots, F_{r_{G_i(L),j,X|i}}]$ are cached at the rest participating (active) groups' caches, except the one that desires them (i.e., except at the caches of group i). This means that each receiver of a specific active group holds in its cache all the file segments that have been requested by the users of the other groups that are active at that particular instance.

[0131] In a step 147 designated "Null-out intra-group interference, Decode own subfile using precoding-based decoder", the receiver applies a standard decoder that reflects the employed state-of-art precoder, to receive its desired subfile. By design of the precoder-and-decoder (which for a specific embodiment can be a simple ZF precoder), a message decoded by a receiver does not contain any file segment intended for another receiver of the same group.

[0132] By combining steps (146) and (147), the receiver can remove all other $K\gamma - 1$ unwanted subfiles, and decode its own subfile.

[0133] Step (148) consists of repeating the process (143) to (147), until the entire desired file is received.

EXAMPLE BASED ON THIRD EMBODIMENT

[0134] Assume a library of $N = 20$ files (for example $N = 20$ different movies) and let there be $K = 20$ wireless users labeled by $1, 2, \dots, 20$. Let each receiving node $k = 1, 2, \dots, 20$ have a cache (storage) of size corresponding to $M = 8$

(equivalent to the size of 8 files), thus corresponding to a normalized cache size of $\gamma = \frac{M}{N} = \frac{8}{20} = \frac{2}{5}$. In theory, the number of users potentially served simultaneously is 8, but state-of-art methods require a file to be segmented in ~ 1 million subfiles. If the constraint on the number of subfiles was 300 then there could only be multicasting gains of 2. To tackle this, and avoid as much of the overhead associated with user cooperation, we will divide users in 5 groups.

3. Phase 1: Placement (first time slot)

- *User enumeration and grouping:* Users are divided into 5 groups of size $L = 4$ users per group. The groups are $G_1 = \{1, 6, 11, 16\}$ (group 1), $G_2 = \{2, 7, 12, 17\}$ (group 2), $G_3 = \{3, 8, 13, 18\}$ (group 3), $G_4 = \{4, 9, 14, 19\}$ (group 4), and $G_5 = \{5, 10, 15, 20\}$ (group 5).

- *Subpacketization indices:* We create the set T that includes $S^* = \binom{\frac{K}{L}}{\frac{K\gamma}{L}} = \binom{5}{2} = 10$ different sets τ ,

where each set τ consists of $Q = \frac{K\gamma}{L} = 2$ distinct numbers from $\{1, 2, \dots, 5\}$ and the set $P = \{1, 2\}$. Thus, $T = \{(1, 2), (1, 3), (1, 4), (1, 5), (2, 3), (2, 4), (2, 5), (3, 4), (3, 5), (4, 5)\}$ and a subfile is defined by three indices, the file associated with it, one number from set P and a set (pair of numbers) from set T .

- *File segmentation:* Each file F_n ($n = 1, 2, \dots, N$) is segmented into $S^* = 20$ subfiles $F_{n,p,\tau}$ where each set τ in the above T is used and $p \in \{1, 2\}$. For example, file F_1 is segmented into subfiles $\{F_{1,1,(1,2)}, F_{1,2,(1,2)}, F_{1,1,(1,3)}, F_{1,2,(1,3)}, F_{1,1,(1,4)}, \dots, F_{1,2,(4,5)}\}$ and similarly with the rest of the files F_2, \dots, F_{20} .

- *First phase transmission:* A base station (B) sequentially transmits x_1 which maps information from F_1, F_2, \dots, F_N .
- *Group-based caching:* Based on x_1 , cache-placement commences, so that each user belonging to group i ,

$(i = 1, 2, \dots, \frac{K}{L})$, caches every subfile $F_{n,p,\tau}$ as long as i can be found inside the set τ . The cache placement indices are defined as follows:

$$Z_1 = \left\{ \begin{array}{l} (1, (1,2)), (2, (1,2)), (1, (1,3)), (2, (1,3)) \\ (1, (1,4)), (2, (1,4)), (1, (1,5)), (2, (1,5)) \end{array} \right\} \text{ (cache for group 1),}$$

$$Z_2 = \left\{ \begin{array}{l} (1, (1,2)), (2, (1,2)), (1, (2,3)), (2, (2,3)) \\ (1, (2,4)), (2, (2,4)), (1, (2,5)), (2, (2,5)) \end{array} \right\} \text{ (cache for group 2),}$$

$$Z_3 = \left\{ \begin{array}{l} (1, (1,3)), (2, (1,3)), (1, (2,3)), (2, (2,3)) \\ (1, (3,4)), (2, (3,4)), (1, (3,5)), (2, (3,5)) \end{array} \right\} \text{ (cache for group 3),}$$

$$Z_4 = \left\{ \begin{array}{l} (1, (1,4)), (2, (1,4)), (1, (2,4)), (2, (2,4)) \\ (1, (3,4)), (2, (3,4)), (1, (4,5)), (2, (4,5)) \end{array} \right\} \text{ (cache for group 4),}$$

$$Z_5 = \left\{ \begin{array}{l} (1, (1,5)), (2, (1,5)), (1, (2,5)), (2, (2,5)) \\ (1, (3,5)), (2, (3,5)), (1, (4,5)), (2, (4,5)) \end{array} \right\} \text{ (cache for group 5).}$$

4. Phase 2: Delivery (second time slot)

- *File requests:* Each user asks for any one file from the library, so user 1 asks for file F_{r_1} , user 2 for F_{r_2} , and so on, where $r_k \in \{1, 2, \dots, 20\}$.
- *First, Group 1 assumes the role of a distributed transmitter assigned with the task to send to all other groups.* It does so by forming all sets of size two from the set $[K] \setminus \{1\} = [5] \setminus \{1\} = \{2, 3, 4, 5\}$, i.e.

$$\chi \in \left\{ \begin{array}{l} (2,3), (2,4), (2,5), \\ (3,4), (3,5), (4,5) \end{array} \right\}.$$

- *Feedback to transmitters and Transmission to groups:* The first transmission is denoted by $x_{1,23}$, where the first index implies that all users belonging in group 1 will transmit messages and the second index implies that all nodes in groups 2 and 3 will be receiving a message. Each active receiving node N_k (i.e., any user from groups 2 or 3) communicates channel state information h_k back to the transmitters. The transmitted vector, from group 1 to the active receivers, then takes the form

$$x_{1,23} = (H_{G_2})^{-1} \cdot \begin{pmatrix} F_{r_2,(1,1,3)} \\ F_{r_7,(1,1,3)} \\ F_{r_{12},(1,1,3)} \\ F_{r_{17},(1,1,3)} \end{pmatrix} + (H_{G_3})^{-1} \cdot \begin{pmatrix} F_{r_3,(1,1,2)} \\ F_{r_8,(1,1,2)} \\ F_{r_{13},(1,1,2)} \\ F_{r_{18},(1,1,2)} \end{pmatrix}.$$

[0135] Then, the received messages - focusing on the case of group 2 ($G_1 = \{2, 7, 12, 17\}$) - stacked as a vector to incorporate the received signals at users $\{2, 7, 12, 17\}$, takes the form:

$$y_{1,23}(G_2) = \begin{pmatrix} F_{r_2,(1,1,3)} \\ F_{r_7,(1,1,3)} \\ F_{r_{12},(1,1,3)} \\ F_{r_{17},(1,1,3)} \end{pmatrix} + H_{G_2} \cdot (H_{G_3})^{-1} \cdot \begin{pmatrix} F_{r_3,(1,1,2)} \\ F_{r_8,(1,1,2)} \\ F_{r_{13},(1,1,2)} \\ F_{r_{18},(1,1,2)} \end{pmatrix}.$$

[0136] As seen from the above equation, each user in Group 2 receives the desired subfile plus a linear combination of all the subfiles intended for the users in the other participating group (group 3 in this case), where this linear combination depends on the channels of these $K\gamma$ participating (active) users. By using the cached content and the knowledge of the channel coefficients, each user can decode their desired subfile.

[0137] In the same manner, we proceed with the rest of the transmissions, i.e., with the rest of the transmissions from Group 1 and then move to the transmissions from Group 2, Group 3 and so on.

[0138] From the above algorithm we can see that there are $K\gamma$ users served at a time, while the number of file segments

$$S^* = \frac{K\gamma}{L} \binom{\frac{K}{L}}{\frac{K\gamma}{L}} = \binom{5}{2} = 20.$$

(subpacketization) is

Without our method, the best known subpacketization needed

$$S^* = K\gamma \binom{K}{K\gamma} = 8 \cdot \binom{20}{8} = 8 \times 125970 \approx 10^6.$$

to achieve this same performance would have been

The

complexity reduces from approximately 1 million (distinct transmissions) to just 30 distinct transmissions.

IV. Forth embodiment (Distributed Wireless)

[0139] There will now be described a preferred forth embodiment which shows the invented role given to distributed base stations each with one or more antennas employed to serve receiving nodes. The new approach combines multi-antenna precoding with cache-aided multicasting and manages to reduce the complexity associated with cache-aided multicasting algorithms and to allow for much reduced transmission delay when delivering files of any given sizes.

[0140] The preferred embodiment, Figure 14, corresponding to a number of base stations K_T 151, each comprised of 1 or more antennas 156, connected via a connection for cooperation and/or coordination 152 and which could be wireless

or wired or hybrid (wireless and wired) and each having a cache size equal to M_T files from a library of N files **153**, serving K users **154**, each fitted with a cache **155** storage of size equal to the size of M files.

[0141] In this embodiment, there are four basic processes embodying two time slots: i) the process of the first time slot to store files in the transmitters' caches (Figure 15), ii) the process, of the first time slot, that leads to content placement in the caches of the receiving nodes (Figure 16), iii) the process of the second time slot that leads to transmission of content by the base stations (Figure 17), and iv) the process of decoding at the receiving nodes (Figure 18).

[0142] Process that leads to content placement in the caches of the transmitting servers: The first process, Figure 15, starts with a step **191** consisting of server-enumeration where the servers are labeled as $s = 1, 2, \dots, K_T$. The second step **192** designated "Pick server $s: s = 1, 2, \dots, K_T$ " is the placement algorithm for the transmitter-side caches.

For $\gamma_T = \frac{M_T}{N}$, the algorithm places, in a step **193**, any subfile at exactly $K_T \times \gamma_T$ transmitting base stations. These $K_T \times \gamma_T$ transmitting base stations --- for that specific subfile --- can play the previous role of the $L = K_T \times \gamma_T$ transmit antennas of the single base station B. The algorithm consecutively caches entire files F_1, F_2, \dots, F_N into the transmitters, such that the first transmitter caches the first M_T files, the second transmitter the next M_T files, and so on, modulo N . Using $Z_{Tx,s}$ to denote the cache of transmitting server $s = 1, 2, \dots, K_T$, then

$$Z_{Tx,s} = \{ F_{1+(n-1) \bmod N} : n \in \{1 + (s-1)M_T, \dots, M_T \cdot s\} \}.$$

[0143] This guarantees the memory constraints and the aforementioned requirement that each subfile resides in $K_T \times \gamma_T$ transmitters.

[0144] Later, in the transmission of any given subfile, these $K_T \times \gamma_T$ transmitters, with knowledge of the CSIT, can now play the role of the aforementioned $L = K_T \times \gamma_T$ antennas and thus can precode this said subfile using the exact same precoders described before, allowing for simultaneous separation of the $L = K_T \times \gamma_T$ streams within any given

group G_g of $L = K_T \times \gamma_T$ receivers. As before, the aforementioned caching allows for treatment of $\frac{K}{L} \gamma + 1 + 1$ groups at a time, and a treatment of $K_T \times \gamma_T + K \times \gamma \leq K$ users at a time.

[0145] The process is terminated in a step **194**, when all caches are full ($s = K_T$) and when each file appears in $L = K_T \times \gamma_T$ base stations.

Process that leads to content placement in the caches of the receiving nodes:

[0146] As illustrated in figure 16, the second process comprises of the following steps **161-166** in sequence to be performed during the first time-slot.

[0147] Steps **161** and **162** are respectively "user enumeration" and "evaluation of grouping function and user grouping".

With respect to these steps, users are enumerated and assigned into one of the $\frac{K}{L}$ groups. The grouping function will take as input, among other things, the user enumeration, the network topology (link strengths), as well as the cache capabilities of the different receiving nodes. In this specific embodiment here, where the link strengths are the same for each user and where the cache sizes are also the same across the different receivers, the grouping function is simple

and it assigns each group with L users, as follows: users $\left\{ 1, \frac{K}{L} + 1, \frac{2K}{L} + 1, \dots, \frac{(L-1)K}{L} + 1 \right\}$ are placed in group

1, users $\left\{ 2, \frac{K}{L} + 2, \frac{2K}{L} + 2, \dots, \frac{(L-1)K}{L} + 2 \right\}$ are placed in group 2 and so on.

[0148] The next step **163** is "Creation of subpacketization indices" which results in the calculation of subpacketization indices τ . To do this, we gather a set T that is some subset of S^* sets τ , where each set τ consists of Q distinct numbers

from $\left\{ 1, 2, \dots, \frac{K}{L} \right\}$. In some specific embodiments, Q can take the value $Q = \frac{K\gamma}{L}$.

[0149] A next step **164**, so-called "File segmentation into subpackets", consists in file segmentation, wherein each file F_n , ($n = 1, 2, \dots, N$) is segmented into S^* subfiles $F_{n,\tau}$, where each set τ in T is used to label a different segment of the file.

[0150] In the next, so called "Transmission of file library", step **165**, the base stations sequentially transmit x_1 , which

maps information from F_1, F_2, \dots, F_{N_r} and each receiving node N_k ($k = 1, 2, \dots, K$) processes their received signal $f_k(x_1)$ to get an estimate of different parts of the library F_1, F_2, \dots, F_{N_r} . The transmission is multicast or partial unicast and it can employ one or more transmit antennas.

[0151] At the completion of step 165, the first process of figure 16 then proceeds to a step 166 designated "Group-based content cache placement", which consists in group-based caching, where each user, depending at least on its group, caches so that the cache overlap - corresponding to the amount of common subfiles - between any two caches, is calculated to be at least a function of whether the two caches belong in the same group or not. In this particular embodiment where the cache sizes are of the same size, and where the link strengths are statistically similar, the caching

is as follows: Each user belonging to group i , (where $i = 1, 2, \dots, \frac{K}{L}$), caches every subfile $F_{n,\tau}$ as long as i can be found inside the set τ .

Process that leads to transmission of content by the base stations:

[0152] With respect to figure 17, there will now be described the third process to be performed in the second time slot, consisting of a sequence of steps 171-179, for gathering file requests (step 171), then creating the transmission indices for active groups (step 172), then identifying K_g active groups (step 173), then gathering channel estimates for these K_g active groups (step 174), then creating the precoding vector for an active group (step 175), repeating for each active group (step 176), then combining the precoding vectors of the K_g active groups (step 177), then transmitting the composite vector simultaneously to all K_g active groups (step 178) and repeating the process for a new set of active groups.

[0153] More particularly, the process of figure 17 starts with a first step 171 being the gathering of file requests, where the receiving nodes request from the transmitters to deliver specific files from the library. In this fourth preferred embodiment, each receiving node N_k requests one file, by communicating to the base stations the information that describes the label of its requested file F_{r_k} ($r_k \in \{1, 2, \dots, N\}$, $k \in K$).

[0154] Then, the process proceeds to a step 172 consisting of gathering transmission indices. In this step, indices χ are created such that to describe which $Q + 1$ groups will be active, which in turn describes to which $Q + 1$ groups to transmit to at any given time. To do this, the process gathers a set X that is some subset of possible sets χ , where each

set χ consists of $Q + 1$ distinct numbers from $\{1, 2, \dots, \frac{K}{L}\}$.

[0155] Then, after completion of step 172, the process proceeds to step 173, which is about identifying the active groups. Here, sets χ are sequentially picked from X . Once a set χ is picked, this immediately identifies the $Q + 1$ active groups, which include group i as long as i can be found in the set χ . Thus each χ automatically identifies the $L \times (Q + 1)$

so called "active receiving nodes". In some specific embodiments, $Q + 1$ can take the value $Q + 1 = \frac{K\gamma}{L} + 1$. Given a χ , the step identifies $(Q + 1) \times L$ active receiving nodes.

[0156] Then, in a step 174, the process performs the gathering of channel estimates for the K_g active groups, where each active receiving node N_k communicates their channel state information h_k to the transmitters. As known by a skilled man, this feedback can also take other forms that include receiver state, precoder preference, location information, etc.

[0157] In step 175, the process creates a precoding vector for an active group, by precoding L subfiles for L users in the active group. Here, for a given active group i , the transmitters that have cached each of the requested files requested by Group i cooperate to generate the L -length signal vector

$$w_{G_i}(\chi) \cdot (H_{G_i})^\perp$$

where $(H_{G_i})^\perp$ is the precoding matrix relative to the channel of group i , and where

$$w_{G_i}(\chi) = [F_{r_{G_i(1),\chi i}}, F_{r_{G_i(2),\chi i}}, \dots, F_{r_{G_i(L),\chi i}}]$$

is the vector comprised of the vector representation of the subfiles $F_{r_{G_i(1),\chi i}}, F_{r_{G_i(2),\chi i}}, \dots, F_{r_{G_i(L),\chi i}}$ wanted by the L users of group i .

[0158] Then, in a step 176 the process repeats the above step (175) for each active group.

[0159] In a step 177 the process combines the precoding vectors $w_{G_i(\chi)} \cdot (H_{G_i})^\perp$ of the K_g active groups, to form a composite vector

5

$$x_{t_s} = \sum_{i \in \chi} w_{G_i(\chi)} \cdot (H_{G_i})^\perp,$$

10 where the summation is over integer numbers i that comprise the currently chosen set χ of active groups.

[0160] In a step 178, the base-stations transmit the composite vector x_{t_s} of length L , simultaneously to all K_g active groups.

[0161] At last, in a step 179, the process repeats the steps 172 -178 for different sets of K_g active groups.

15 **Process of decoding at the receiving nodes:**

[0162] With respect to figure 18, there will now be described the third process- which is presented here from the point of view of any one user for decoding at the receiving nodes and which comprises the sequence of steps 181-188 in sequence.

20 [0163] The fourth process starts with step 181, wherein a test is performed to determine whether the user belongs to the upcoming transmission, which in turn means that the receiver is active.

[0164] The result of the test is processed in a step 182.

[0165] If the user does not belong to the upcoming transmission, then the process proceeds to step 188 wherein the receiver will wait until the next transmission.

25 [0166] If the user belongs to the upcoming transmission, then the process proceeds to a step 183, the receiver will broadcast its channel state information.

[0167] Then, in a step 184, the receiver will listen for and record the other channel information from the other active users.

[0168] Then, the process proceeds to a step 185, wherein the receiver records the received signal, which is meant for the active groups.

30 [0169] In step 186, the receiver uses channel information and cached content to remove (cache-out) inter-group interference. This is achieved by using a state-of-art cache-aided multicasting method well known to a skilled man, and which is elevated in this step to a group-level such that, by design, the sets T and X , guarantee that all the elements of

35 $w_{G_i(\chi)} = [F_{r_{G_i(1),\chi i}}, F_{r_{G_i(2),\chi i}}, \dots, F_{r_{G_i(L),\chi i}}]$ are cached at the rest participating (active) groups' caches, except the one that desires them (i.e., except at the caches of group i). This means that each receiver of a specific active group holds in its cache all the file segments that have been requested by the users of the other groups that are active at that particular instance.

40 [0170] In a step 187, the receiver applies a standard decoder that reflects the employed state-of-art precoder, to receive its desired subfile. By design of the precoder-and-decoder (which for a specific embodiment can be a simple ZF precoder), a message decoded by a receiver does not contain any file segment intended for another receiver of the same group.

[0171] By combining steps (186) and (187), the receiver can remove all other Q unwanted subfiles and decode its own subfile.

45 [0172] Step (188) consists of repeating the process (181) to (187), until the entire desired file is received.

EXAMPLE BASED ON FOURTH EMBODIMENT

[0173] Assume 8 base stations with 1 transmit antenna each, being deployed in a mall. Each transmitter has access

50

to half the library, i.e.: $\gamma_T = \frac{1}{2}$, thus $L = S \times K\gamma = 4$. Let there be a library of $N = 20$ files (for example $N = 20$ different movies) and let there be $K = 20$ wireless users (previously referred to as receiving nodes) labeled by 1,2,...,20. Let each receiving node $k = 1, 2, \dots, 20$ have a cache (storage) of size corresponding to $M = 8$ (equivalent to the size of 8 files),

55

thus corresponding to a normalized cache size of $\gamma = \frac{M}{N} = \frac{8}{20} = \frac{2}{5}$.

◦ Phase 1: Transmitter Side Placement (first time slot)

a. Transmitters 1, 3, 5, 7 cache files $\{F_1, F_2, \dots, F_{10}\}$ and Transmitters 2,4,6,8 cache files $\{F_{11}, F_{12}, \dots, F_{20}\}$.

5 ◦ Phase 2: Receiver Side Placement (first time slot)

▪ *User enumeration and grouping:* Users are divided into 5 groups of size $L = 4$. The groups are $G_1 = \{1,6,11,16\}$ (group 1), $G_2 = \{2,7,12,17\}$ (group 2), $G_3 = \{3,8,13,18\}$ (group 3), $G_4 = \{4,9,14,19\}$ (group 4), and $G_5 = \{5,10,15,20\}$ (group 5).

10

$$S^* = \binom{\frac{K}{L}}{\frac{KY}{L}} = \binom{5}{2} = 10 \text{ different sets } \tau,$$

▪ *Subpacketization indices:* We create the set T that includes $Q = \frac{KY}{L} = 2$ where each set τ consists of $Q = \frac{KY}{L} = 2$ distinct numbers from $\{1,2,\dots,5\}$. Thus $T = \{(1,2), (1,3), (1,4), (1,5), (2,3), (2,4), (2,5), (3,4), (3,5), (4,5)\}$.

15

▪ *File segmentation:* Each file F_n ($n = 1,2, \dots, N$) is segmented into $S^* = 10$ subfiles $F_{n,\tau}$ where each set τ in the above T is used. For example, file F_1 is segmented into subfiles $\{F_{1,(1,2)}, F_{1,(1,3)}, F_{1,(1,4)}, \dots, F_{1,(4,5)}\}$ and similarly with the rest of the files F_2, \dots, F_{20} .

20

▪ *First phase transmission:* The base stations sequentially transmit x_1 which maps information from F_1, F_2, \dots, F_N .

▪ *Group-based caching:* Based on x_1 , cache-placement commences, so that each user belonging to group i

25

$\left(i = 1,2, \dots, \frac{K}{L}\right)$ caches every subfile $F_{n,\tau}$ as long as i can be found inside the set τ . The cache placement is defined by the following: $Z_1 = \{(1,2), (1,3), (1,4), (1,5)\}$ (cache for group 1), $Z_2 = \{(1,2), (2,3), (2,4), (2,5)\}$ (cache for group 2), $Z_3 = \{(1,3), (2,3), (3,4), (3,5)\}$ (cache for group 3), $Z_4 = \{(1,4), (2,4), (3,4), (4,5)\}$ (cache for group 4), and $Z_5 = \{(1,5), (2,5), (3,5), (4,5)\}$ (cache for group 5). This means that for example, all users of group 1 (i.e., users 1,6,11,16) will cache all subfiles $F_{n,(1,2)}, F_{n,(1,3)}, F_{n,(1,4)}, F_{n,(1,5)}$, for each file F_n , $n = 1,2, \dots, 20$. Similarly all users of group 2 (i.e., users 2,7,12,17) will cache all subfiles $F_{n,(1,2)}, F_{n,(2,3)}, F_{n,(2,4)}, F_{n,(2,5)}$, again for each file F_n , and so on. It can be seen that inter-group caches have 1 out of 5 parts in common (e.g. all users in groups 1 and 2 have segments labeled by (1,2) in common), while intra-group caches have all elements in common (i.e., user caches of the same group are identical).

35

◦ Phase 2: Delivery (second time slot)

▪ *File requests:* Each user asks for any one file from the library, so user 1 asks for file F_{r_1} , user 2 for F_{r_2} , and so on, where $r_k \in \{1,2,\dots,20\}$.

40

▪ *Transmission indices and active group identification:* We consider the set $X = \{(1,2,3), (1,2,4), (1,2,5), (1,3,4), (1,3,5), (1,4,5), (2,3,4), (2,3,5), (2,4,5), (3,4,5)\}$ consisting of sets χ , where each set χ consists of $Q + 1 = 3$ distinct numbers from the set $\{1,2,3,4,5\}$. Each different χ describes a different set of $Q + 1 = 3$ active groups, corresponding to $(Q + 1) \times L = 12$ active users. For example, when $\chi = (1,2,3)$, it means that the active groups are groups 1,2 and 3, which are the groups that the base stations will transmit to, during this time.

45

▪ *Feedback to transmitters and Transmission to groups:* The first transmission is denoted by x_{123} where the index implies that all users belonging in groups 1,2 and 3 are going to be receiving a desired subfile. Each active receiving node N_k (i.e., any user from groups 1,2 or 3) communicates channel state information h_k back to the transmitter. The transmitted vector, from the base stations to the active receivers, then takes the form

50

$$x_{123} = (H_{G_1})^{-1} \cdot \begin{pmatrix} F_{r_1,(2,3)} \\ F_{r_6,(2,3)} \\ F_{r_{11},(2,3)} \\ F_{r_{16},(2,3)} \end{pmatrix} + (H_{G_2})^{-1} \cdot \begin{pmatrix} F_{r_2,(1,3)} \\ F_{r_7,(1,3)} \\ F_{r_{12},(1,3)} \\ F_{r_{17},(1,3)} \end{pmatrix} + (H_{G_3})^{-1} \cdot \begin{pmatrix} F_{r_3,(1,2)} \\ F_{r_8,(1,2)} \\ F_{r_{13},(1,2)} \\ F_{r_{18},(1,2)} \end{pmatrix}.$$

55

[0174] Then, the received messages - focusing on the case of group 1 ($G_1 = \{1,6,11,16\}$) - stacked as a vector to incorporate the received signals at users $\{1,6,11,16\}$, takes the form:

$$y_{123}(G_1) = \begin{pmatrix} F_{r_1,(2,3)} \\ F_{r_6,(2,3)} \\ F_{r_{11},(2,3)} \\ F_{r_{16},(2,3)} \end{pmatrix} + H_{G_1} \cdot (H_{G_2})^{-1} \cdot \begin{pmatrix} F_{r_2,(1,3)} \\ F_{r_7,(1,3)} \\ F_{r_{12},(1,3)} \\ F_{r_{17},(1,3)} \end{pmatrix} + H_{G_1} \cdot (H_{G_3})^{-1} \cdot \begin{pmatrix} F_{r_3,(1,2)} \\ F_{r_8,(1,2)} \\ F_{r_{13},(1,2)} \\ F_{r_{18},(1,2)} \end{pmatrix}.$$

[0175] As seen from the above equation, each user in Group 1 receives the desired subfile plus a linear combination of all the subfiles intended for the users in the other participating groups (groups 2 and 3 in this case), where this linear combination depends on the channels of these $K\gamma + L$ participating (active) users. By using the cached content and the knowledge of the channel coefficients, each user can decode their desired subfile.

[0176] In the same manner, we proceed with the rest of the transmissions, i.e., with the rest of the χ , i.e., with $\chi = (1,2,4)$ corresponding to active groups 1,2 and 4, then with $\chi = (1,2,5)$ corresponding to active groups 1,2 and 5, and so on.

[0177] From the above algorithm we can see that there are $K\gamma + L = 12$ users served at a time, while the number of

$$S^* = \binom{\frac{K}{L}}{\frac{K\gamma}{L}} = \binom{5}{2} = 10.$$

file segments (subpacketization) is needed to achieve this same performance would have been

$$S^* = \binom{K}{K\gamma} \times \binom{K_T}{L} = \binom{20}{8} \times \binom{4}{2} = 125970 \times 6 = 755820.$$

[0178] The complexity reduces from approximately 7 million (distinct transmissions), to just 10 distinct transmissions.

Claims

1. A process of communication between a transmitting base station (B) with L transmit antennas, and a set of receiving nodes N_1, N_2, \dots, N_K in a wireless communication system, wherein nodes N_1, N_2, \dots, N_K have storage capabilities and wherein the said process involves the steps of:

- grouping K users into groups, with a specified number of users per group;
- identifying a set of files of interest F_1, F_2, \dots, F_N , that may be requested by the receiving nodes;
- arranging a first time slot (time slot 1) during which said base station (B) transmits a signal x_1 corresponding to the downlink transmission with said receiving nodes, where said signal x_1 maps information from F_1, F_2, \dots, F_N , and where said signal x_1 is received, processed and partially stored (cached) by said receiving nodes, where caching is at least a function of said grouping;
- arranging a second time slot during which some or all receiving nodes N_k ($k = 1,2,\dots,K$) request additional information, each in the form of files F_{rk} , during which some or all of the receivers feed back some information about the channel and during which the base station (B) generates a multi-antenna transmit signal x_2 with the purpose of transmission of requested data files to said receiving nodes;
- executing into some or all said receiving nodes N_k ($k = 1,2,\dots,K$) a decoding process using the information received in said first and second time slots and in the caches, so as to allow each receiving node N_k to compute an estimation of their requested file F_{rk} , whereby the accuracy of the decoding process is improved by prudent utilization of the said steps of grouping, caches, first time slot transmission, and second time-slot transmission.

2. The process according to claim 1 wherein, in the first time slot, segments of files are cached as a function of at least the group that a receiving node belongs to.

3. The process according to claim 1 wherein, in the second time slot, the requested file segments $F_{r1}, F_{r2}, \dots, F_{rk}$ are mapped onto signals transmitted from multiple antennas over a wireless channel, where said mapping considers user grouping information, the resulting placement of files in caches, and available channel state information.

4. The process according to claim 1 **characterized in that** it includes the step of executing a decoding process in

some or all said receiving nodes N_1, N_2, \dots, N_K so as to allow each such node N_k to generate an estimation of file F_{r_k} .

- 5 5. The process according to claim 3 where said group-based mapping considers available receiver-state information, receiver location information and/or network connectivity information, or where the feedback is in the form of selection between predefined precoders.
6. The process according to claim 3 and claim 2 wherein each receiver has at least two antennas and wherein there are more than one transmitters with one or more transmit antennas.
- 10 7. The process according to claim 3 and claim 2 also involving wired communication and wherein the transmitting server or one or more intermediate transmitting nodes are connected directly to at least two receiving nodes.
8. The process according to claim 7 wherein the receiving nodes have at least two shared channels.
- 15 9. The process according to claim 1 wherein some of the receiving nodes may not have storage capabilities (no cache).
10. A transmitting apparatus adapted for carrying out the process of claim 1, further comprising:
 - 20 - a file segmenter implemented as a software or hardware or a combination of the two, where it receives one or more files and divides it into smaller parts;
 - a precoder implemented as a software or a dedicated hardware or some combination of the two, which creates a vector or matrix of signals by combining requested file segments with information fed back from the receivers;
 - a combiner which is implemented as a software or hardware or a combination of the two and which receives precoded message vectors and combines them into one vector;
 - 25 - an encoder which is implemented as a software or hardware or a combination of the two and receives a vector from the combiner and maps it into an encoded message.
 - a multi-antenna transmitter that takes as input the signals from the encoder, and transmits them over the channel.
- 30 11. The apparatus of Claim 10 where the transmitter is a network MIMO device comprised additionally of:
 - a communicator which allows for exchange of information and/or synchronization between transmitters, which could be in the form of wired infrastructure and/or wireless infrastructure.
 - a distributed precoder.
- 35 12. A receiving apparatus adapted for carrying out the process of claim 1, further comprising:
 - a cache; Each receiving node N_k , ($k = 1, 2, \dots, K$) with storage capabilities can store content in their cache, where this content is stored during the said first time slot, during and after the receiving node receives and processes said signal x_1 , and after the receiving node estimates information regarding files F_1, F_2, \dots, F_N ;
 - 40 - a File requester; This is a method that allows each receiving node N_k to communicate to the base station B at the beginning of time slot 2, the information that describes the label of their requested file F_{r_k} , ($r_k \in \{1, 2, \dots, N\}$).
 - a CSI communicator: This is a method that allows the receiving nodes to communicate to the base station B the information that describes their channel state (or location state, or connectivity state, or preferred precoder).
 - 45 - decoder which executes the decoding process in some or all said receiving nodes N_1, N_2, \dots, N_K and which takes as input the processed signal x_2 , as well as the processed signal x_1 including the information that was stored at the cache of the corresponding receiving node.
- 50 13. A process of communication between K independent nodes N_1, N_2, \dots, N_K , where said nodes are requesting at least one out of a plurality of files from a known library and where at least some nodes have caching capabilities and wherein a number of the transmitting nodes are cooperating for simultaneous transmission wherein the process determines caching at the nodes as a function of, among other things, the maximum number of file segments per file.
- 55 14. The process of claim 13 where the caching at the nodes is a function of the number of cooperating nodes, or a function of the group a node belongs to.
15. A process of wireless communication, where there are at least two base stations, each with one or more antennas, communicating to K receiving nodes from some library of files and where each base station can store part of this

library and where each file is stored in at least two base stations, wherein said process is further comprised of:

- a process of file placement at the transmitters, such that each file can be found in at least two base stations;
- a process of creating distributed precoders between the base stations;

5

16. An apparatus comprised of

- a communicator which could be in the form of wireless network or a wired connection or a combination of the two, between the base stations which allows for the exchange of some form of channel state information; and some form of cooperation so as to create distributed precoders;
- a distributed precoder, which takes input from the communicator and designs a signal such that when combined with the others in the wireless medium will take the form of a vector transmitted from a multi-antenna base station.

10

15

20

25

30

35

40

45

50

55

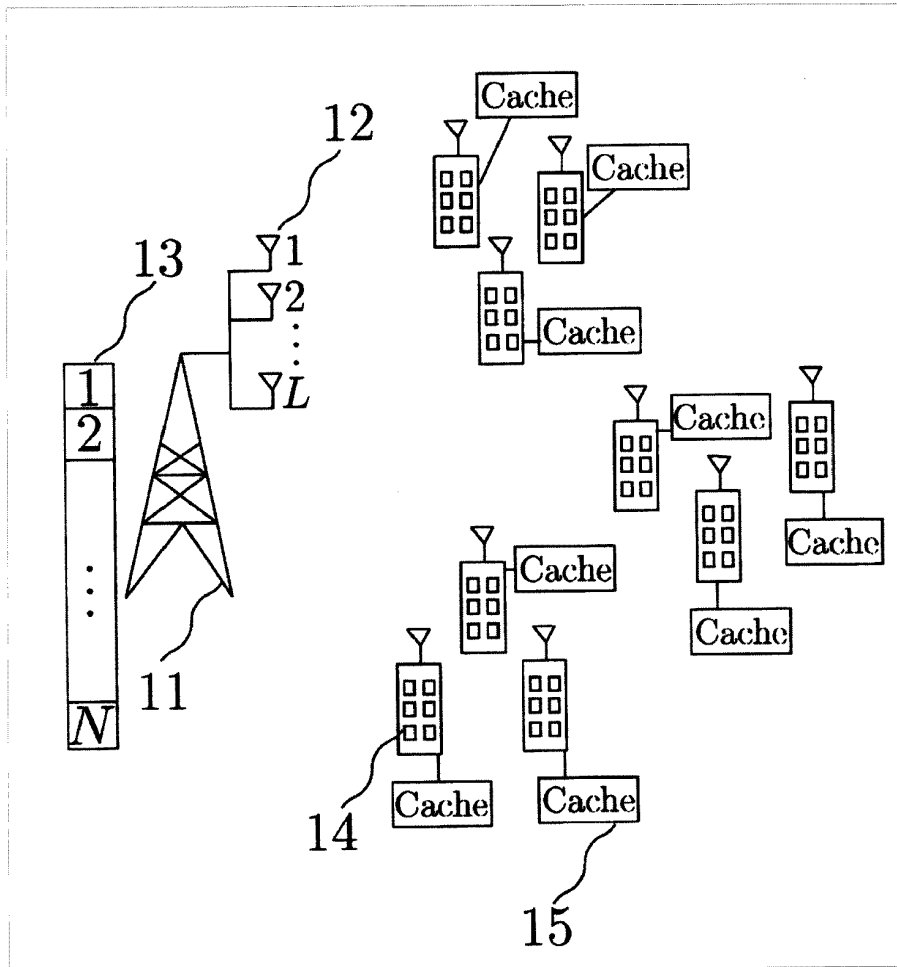


Fig. 1

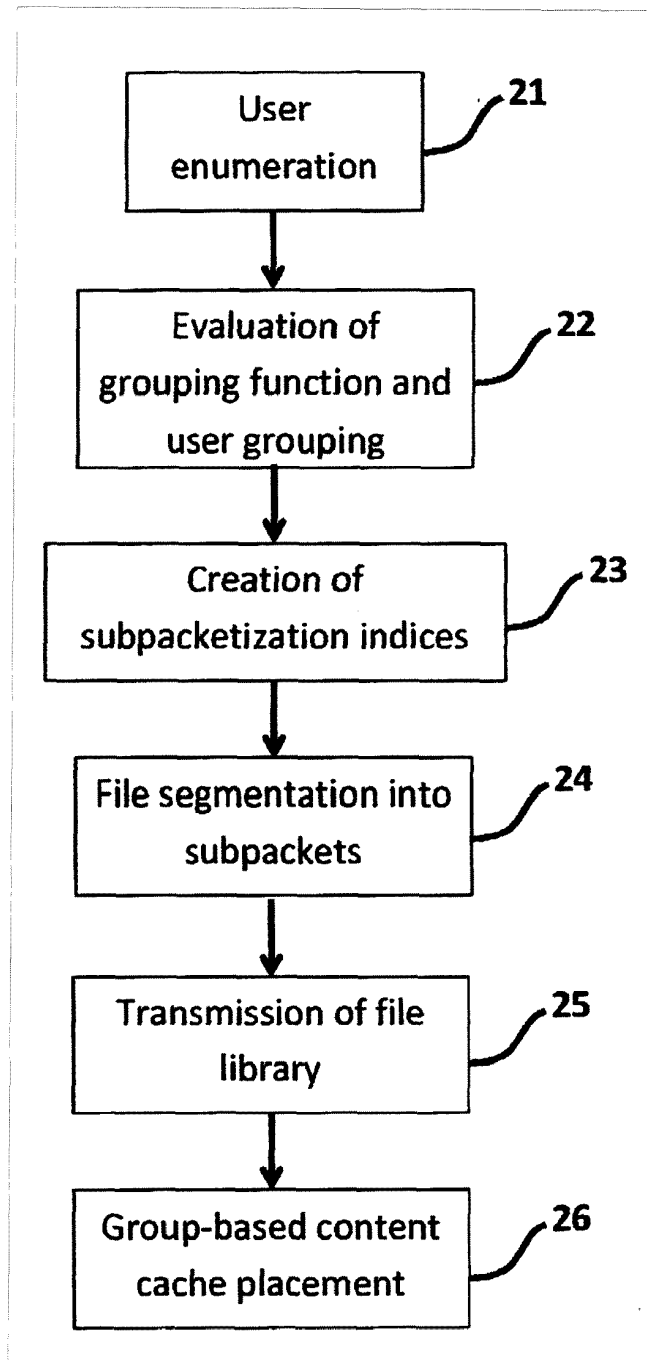


Fig. 2

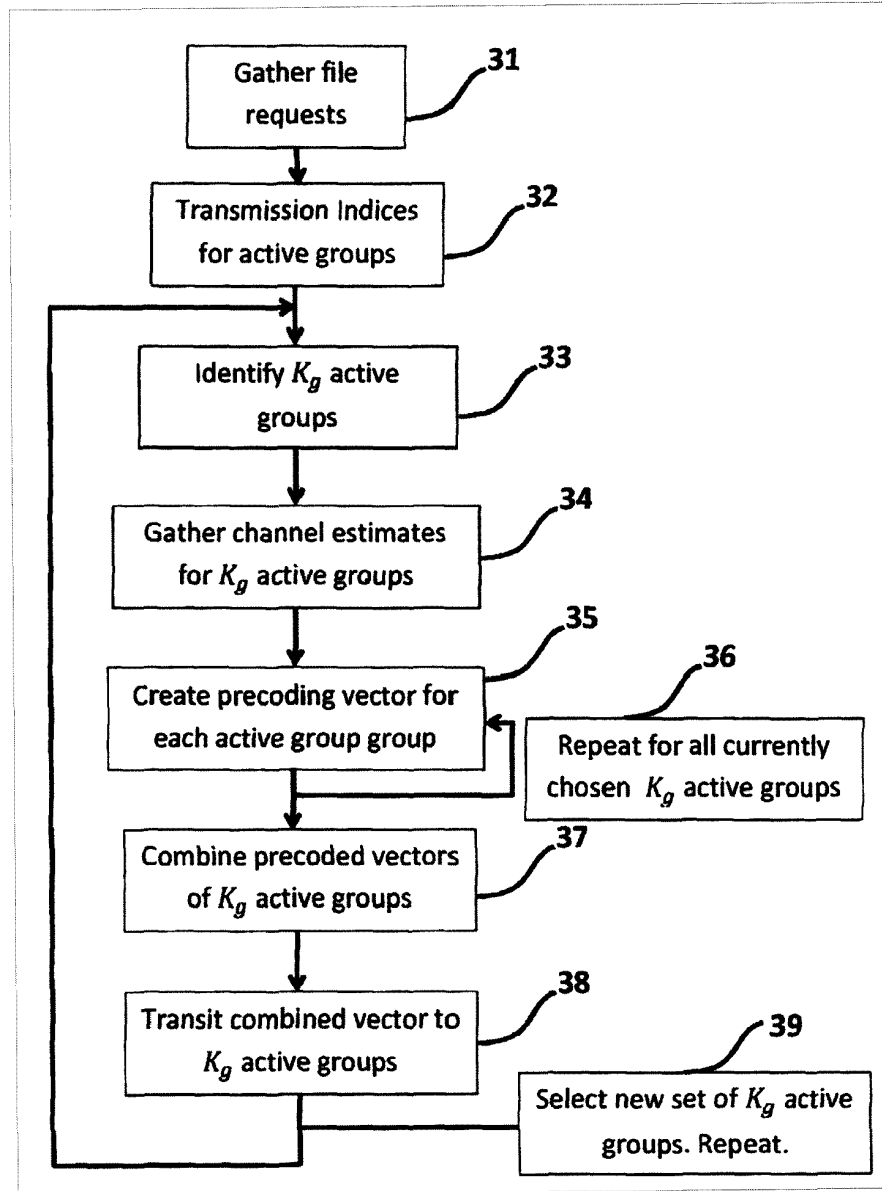


Fig. 3

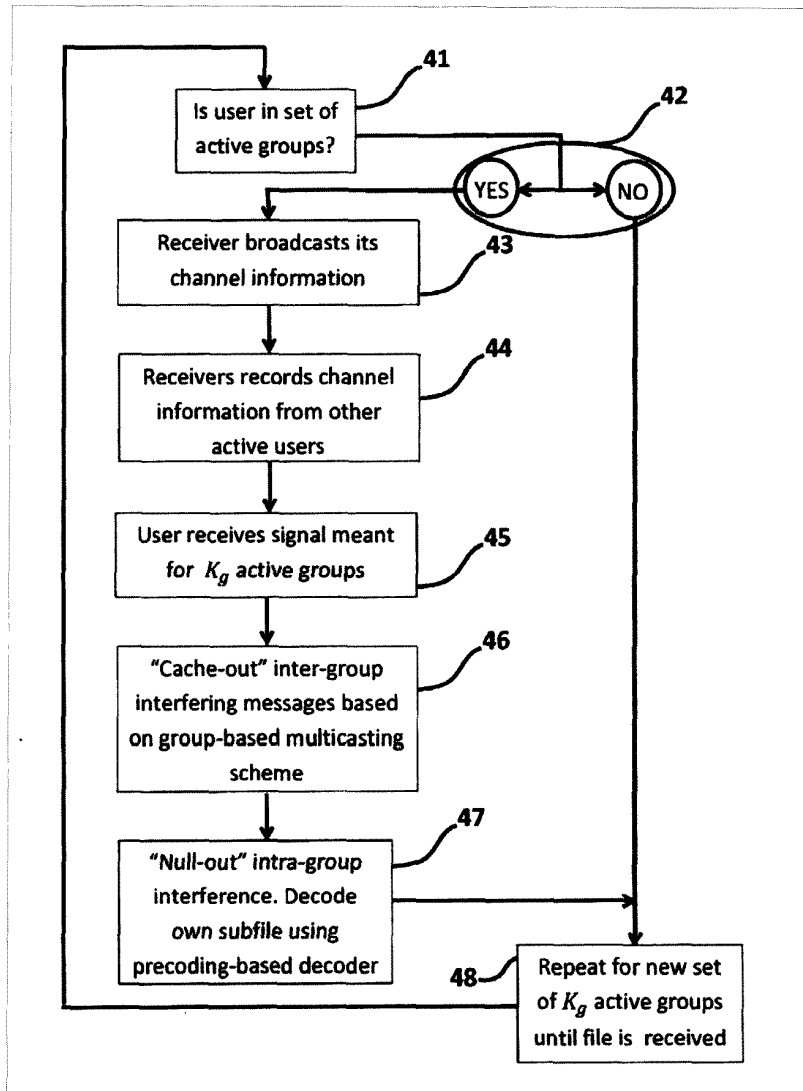


Fig. 4

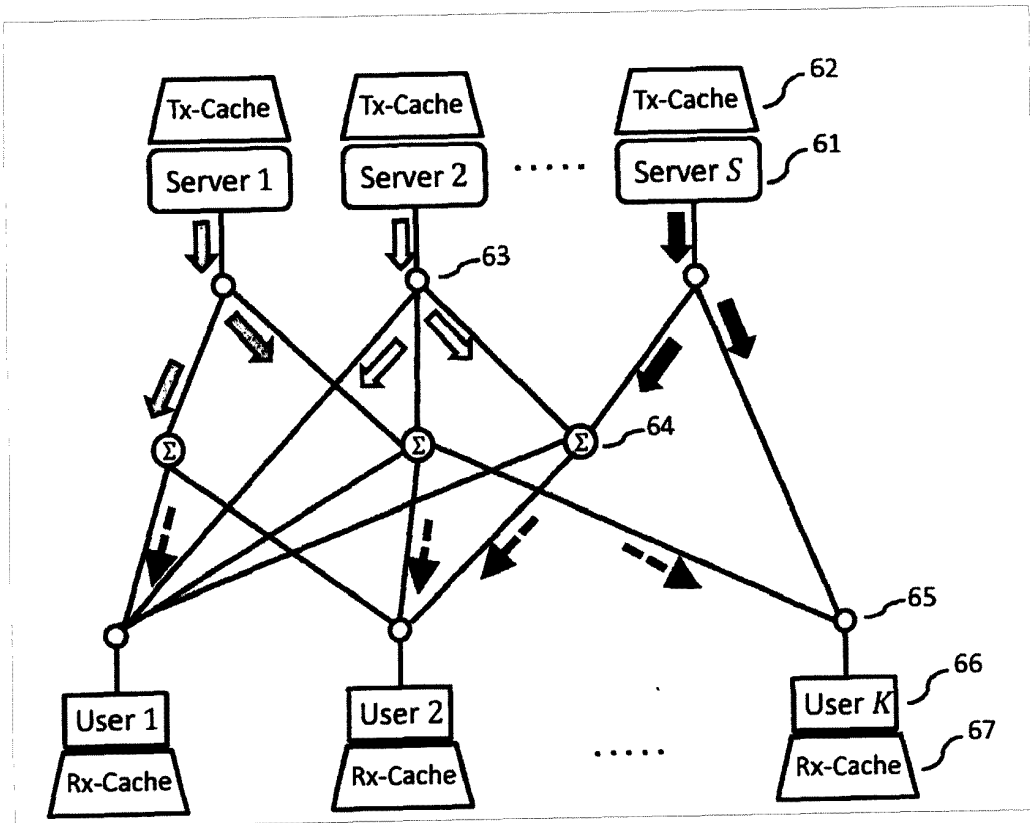


Fig. 5

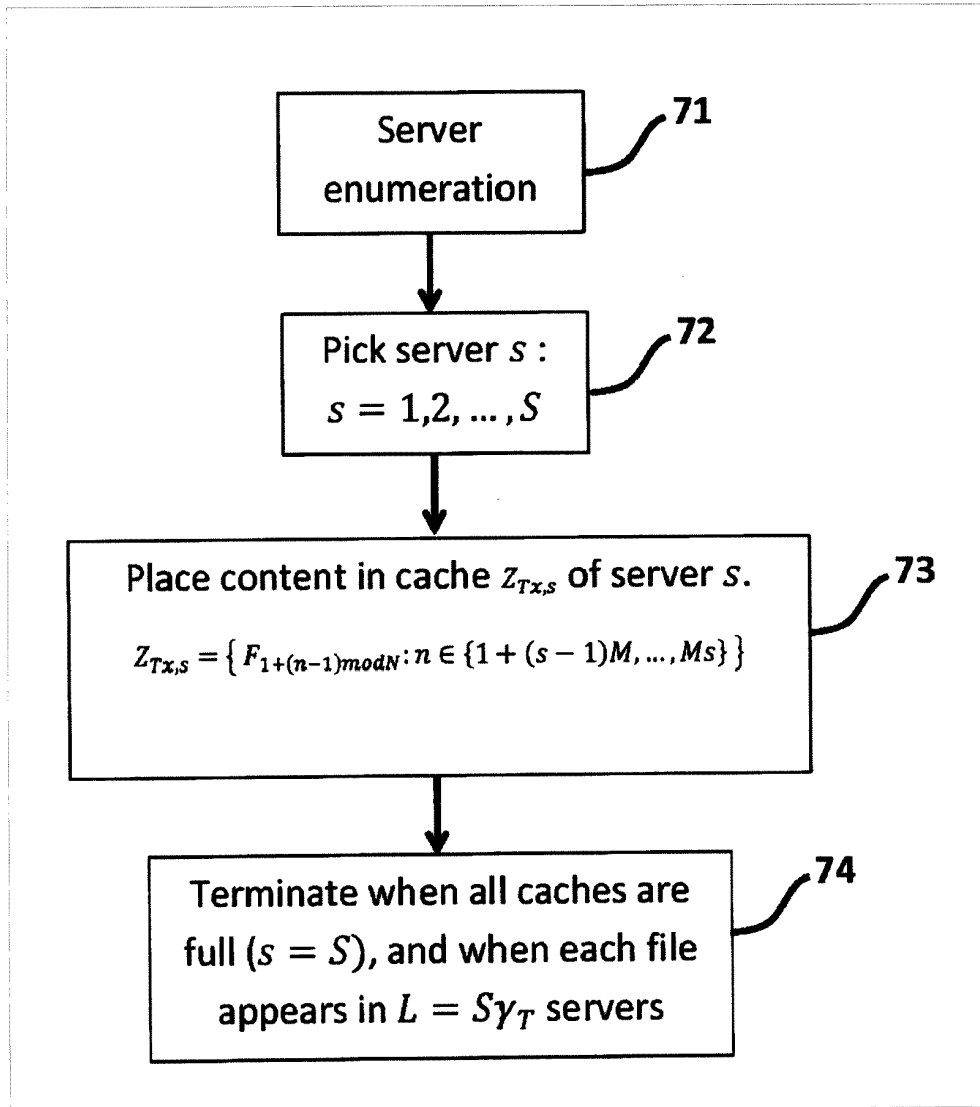


Fig. 6

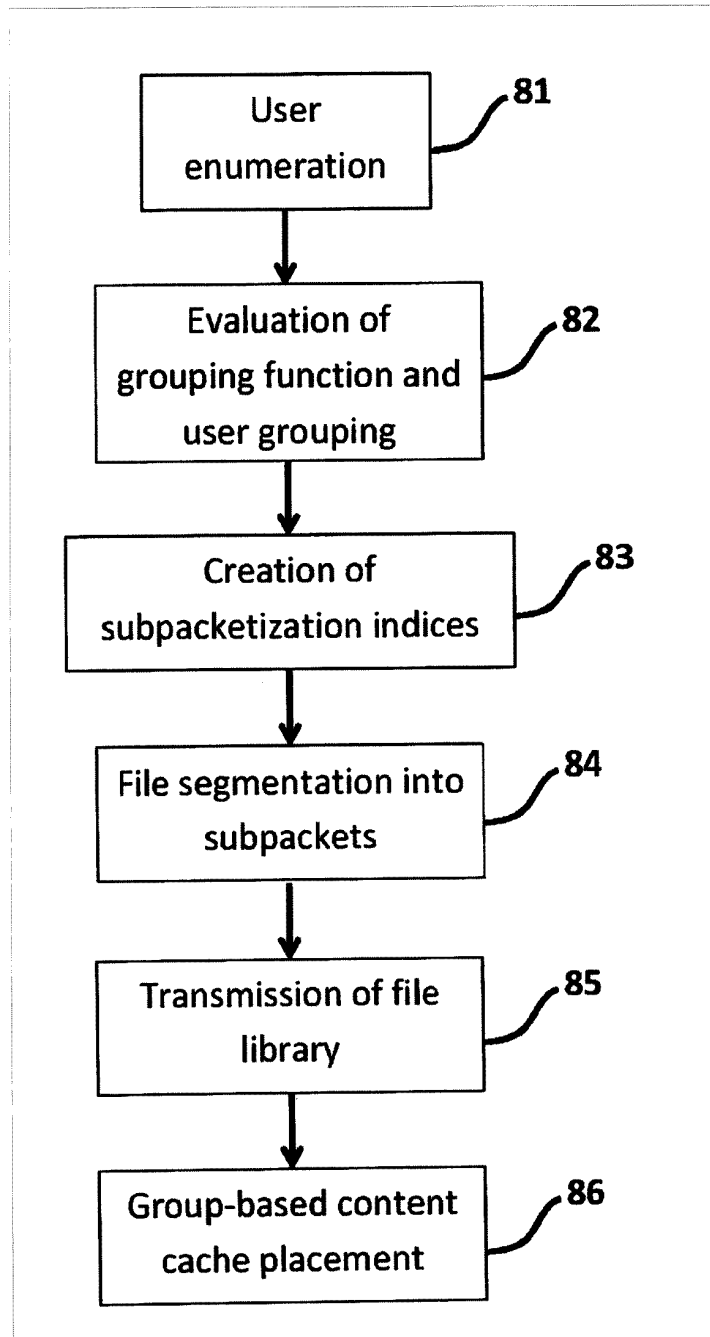


Fig. 7

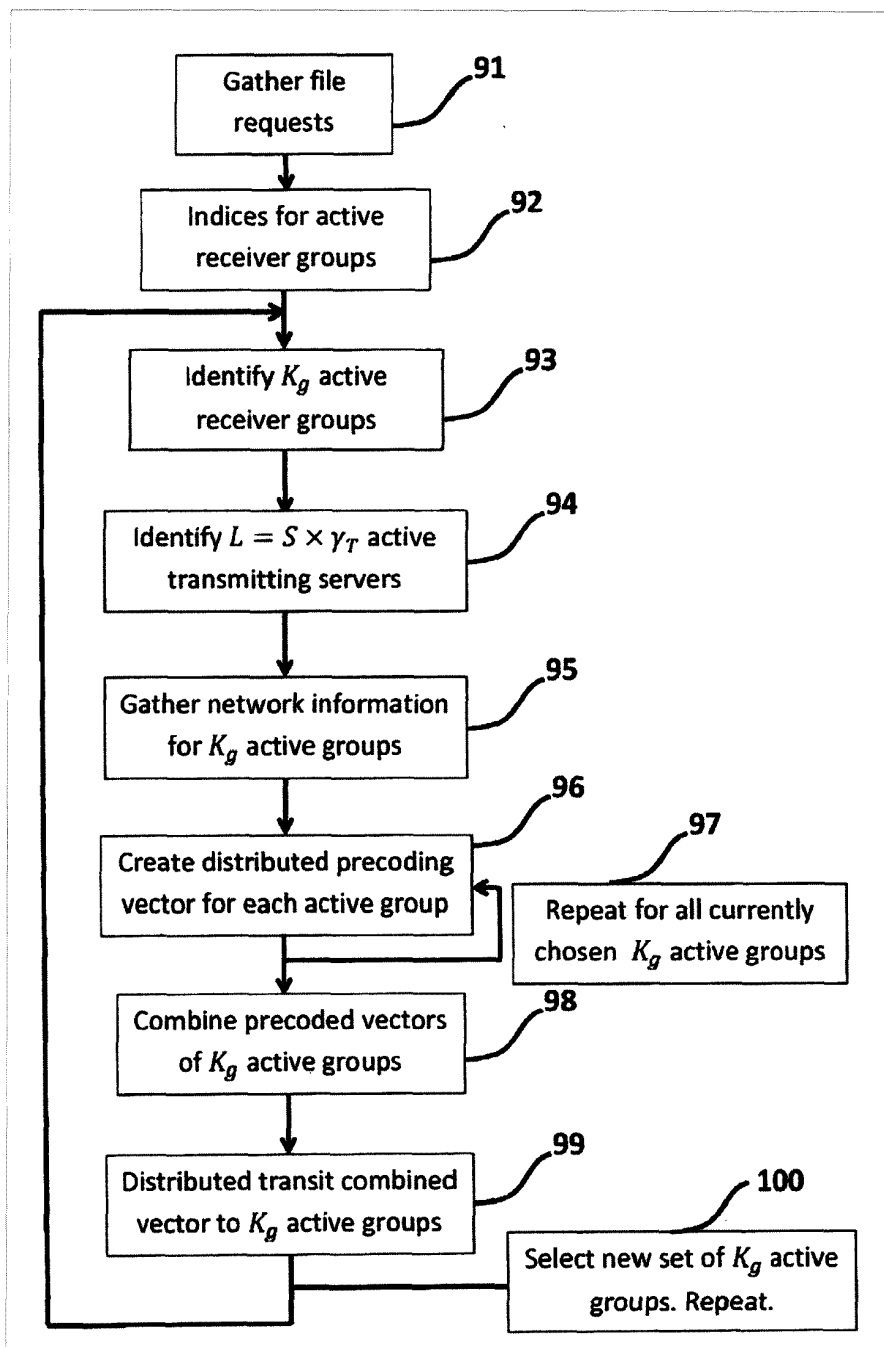


Fig. 8

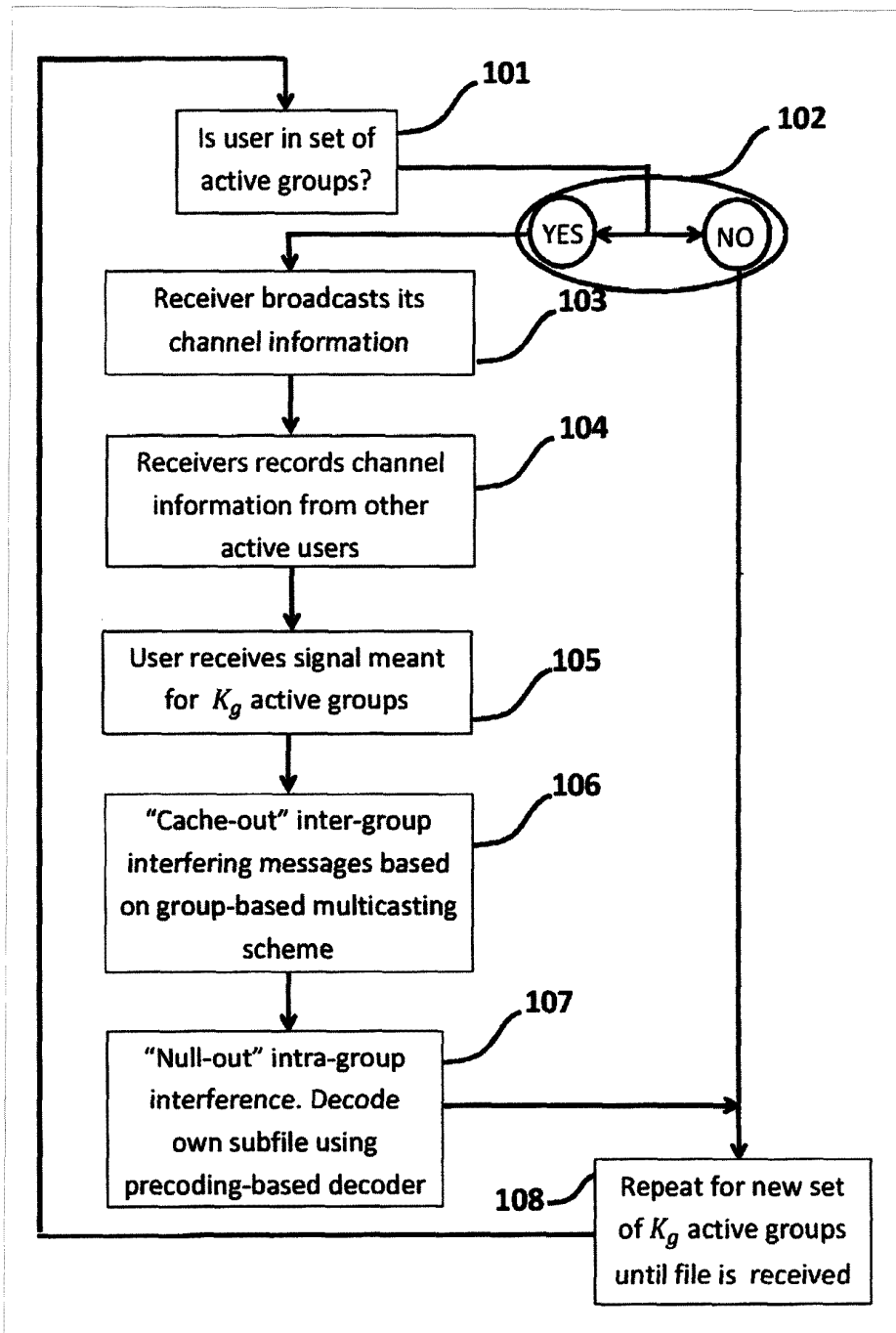


Fig. 9

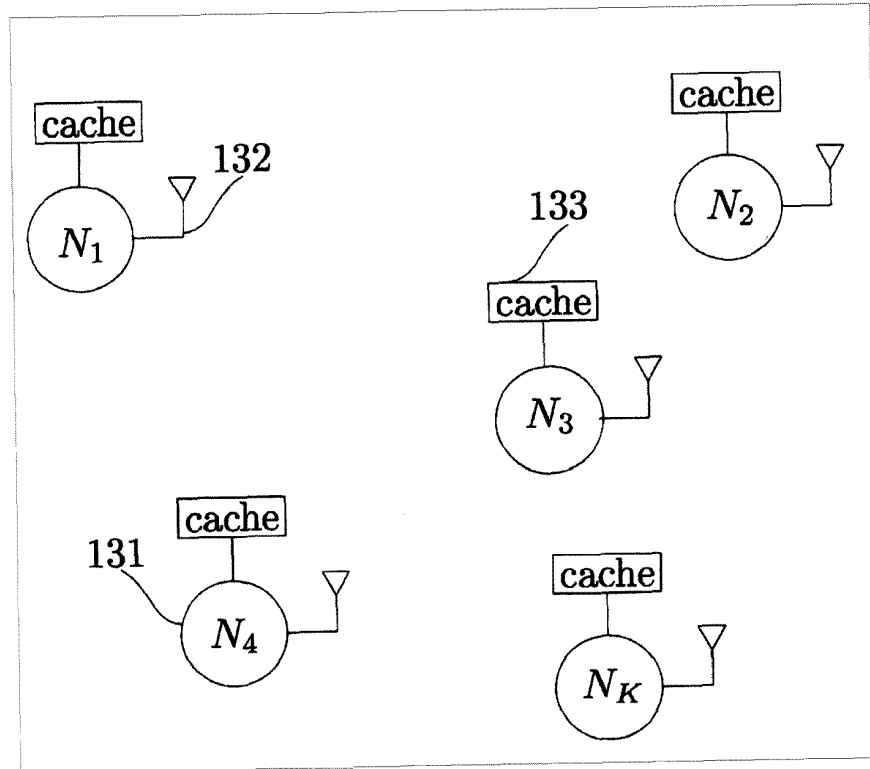


Fig.10

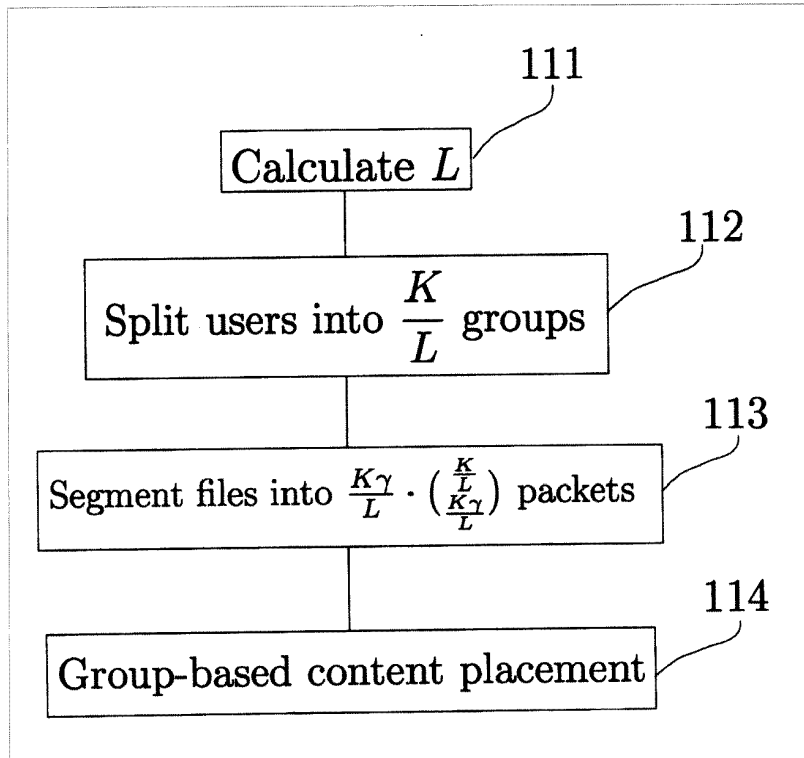


Fig.11

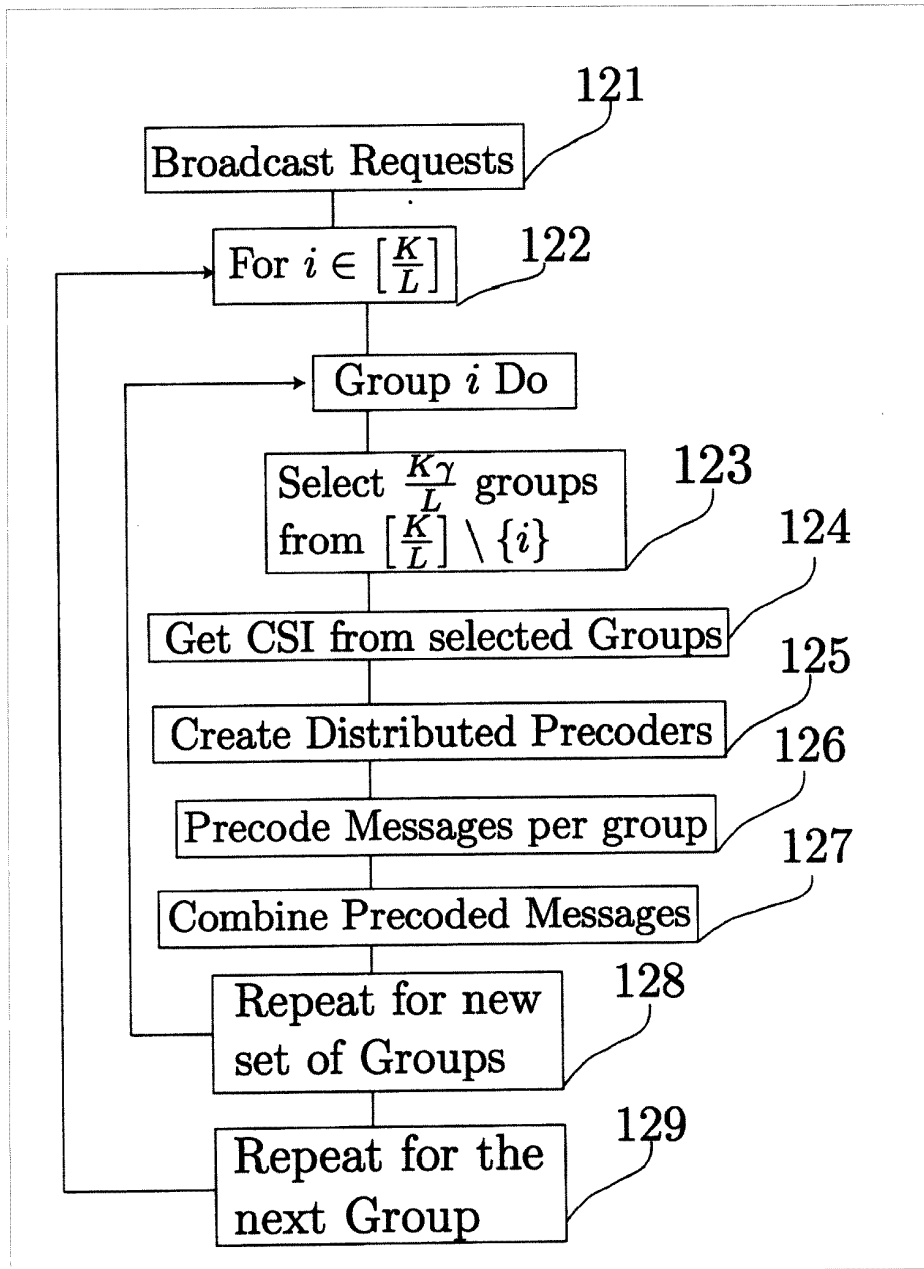


Fig. 12

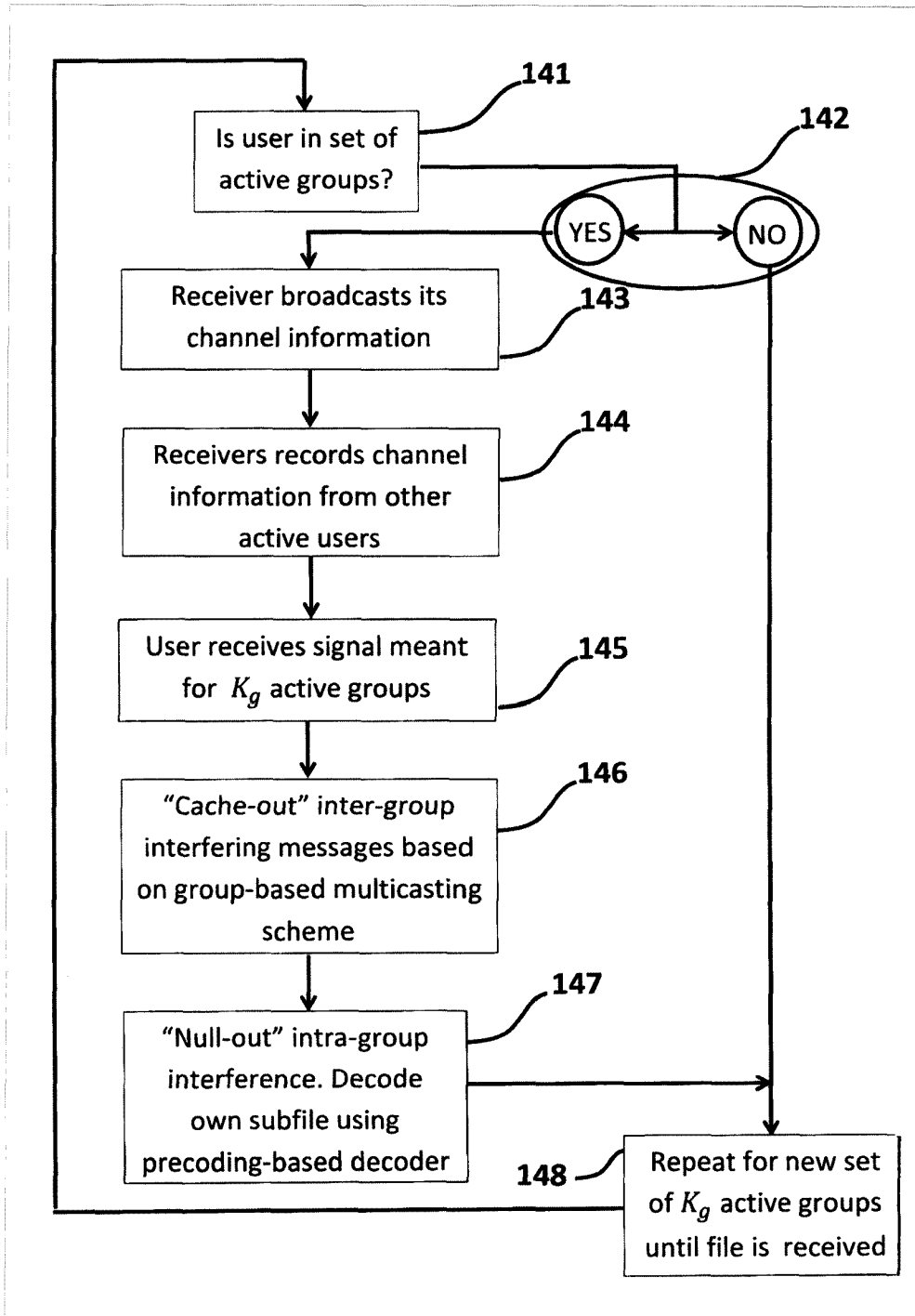


Fig. 13

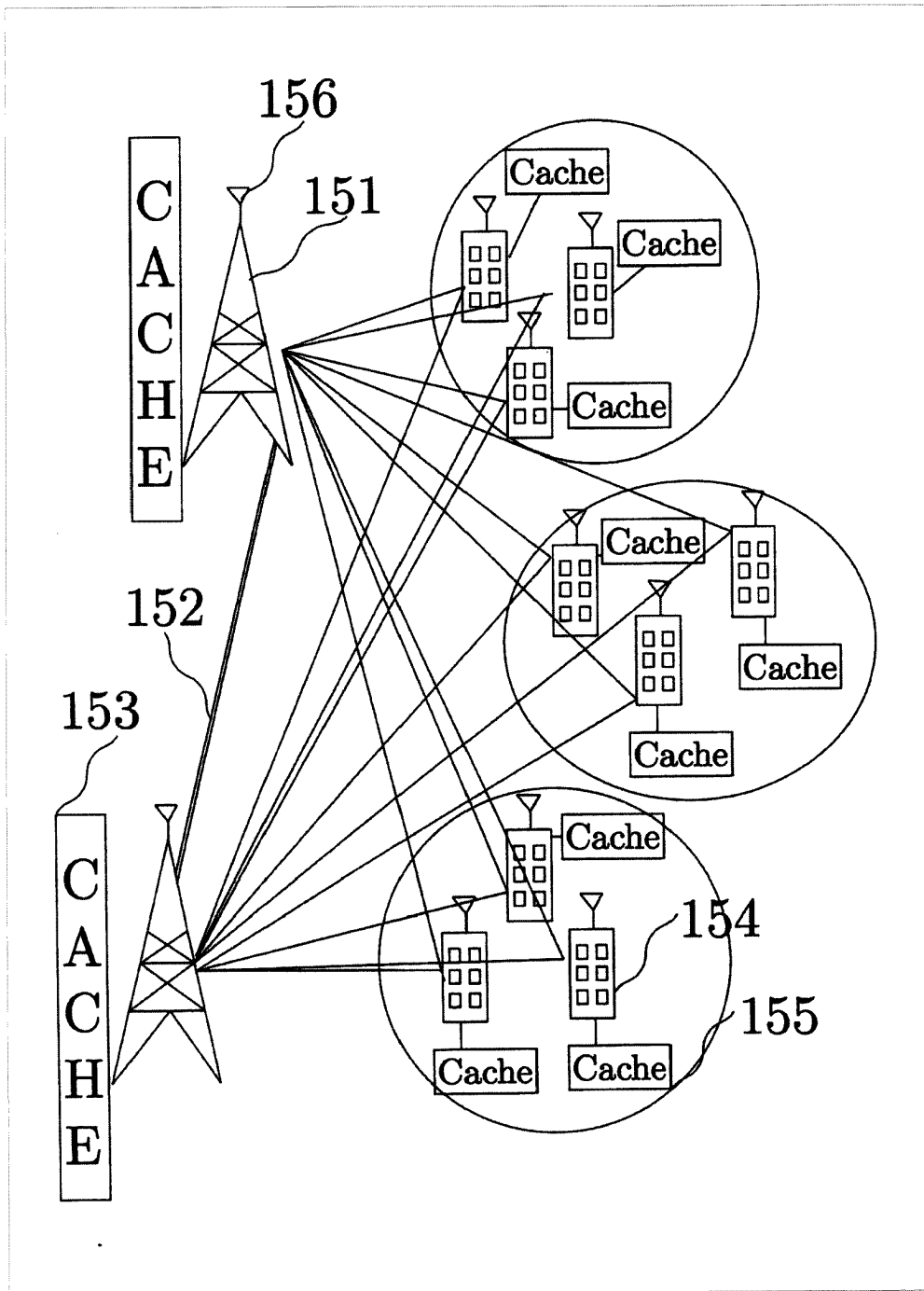


Fig. 14

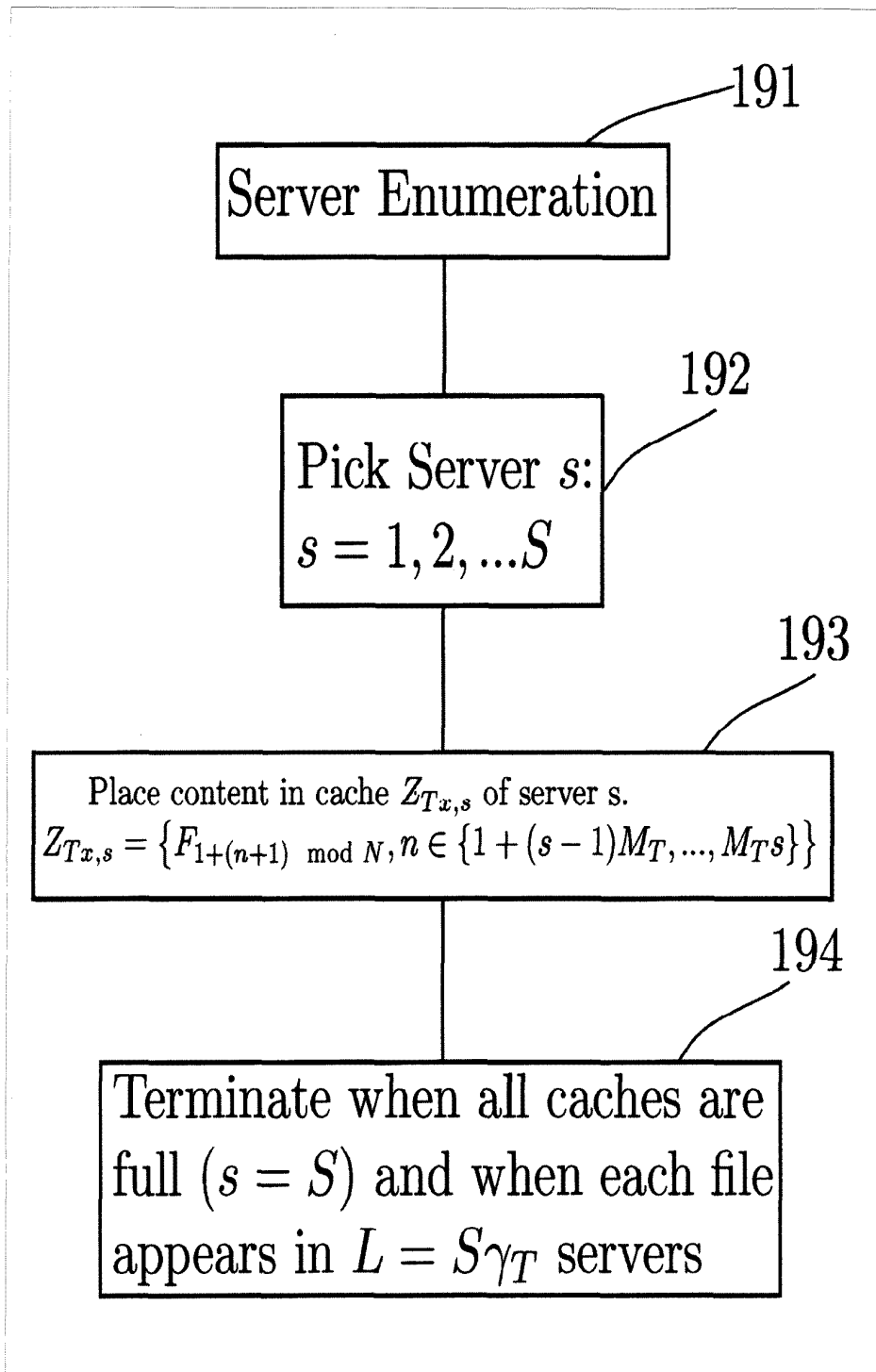


Fig. 15

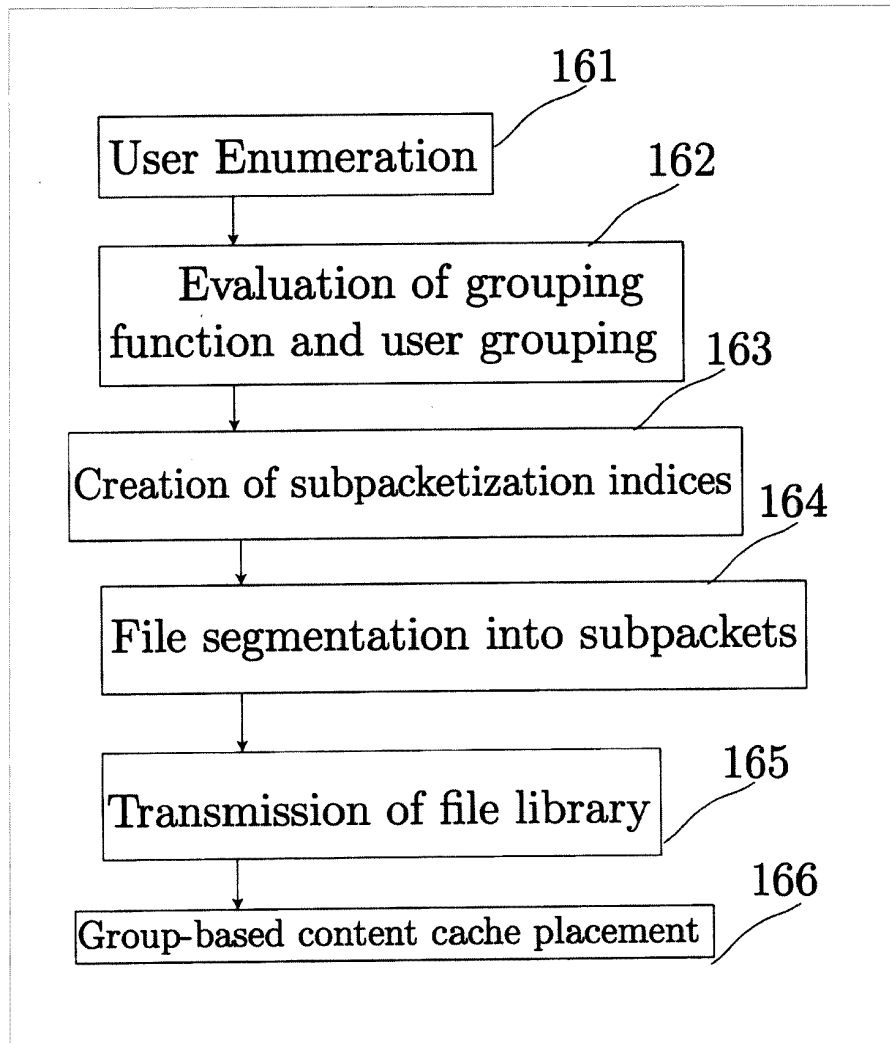


Fig. 16

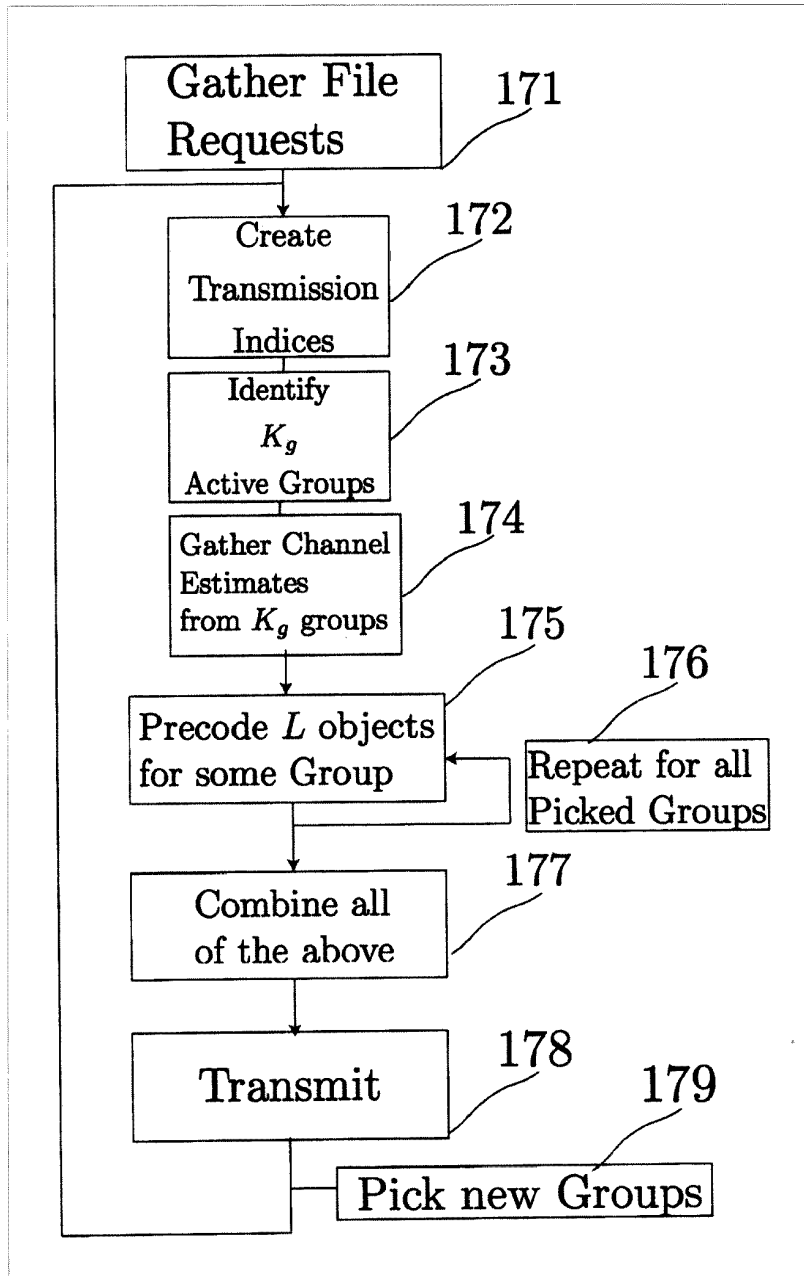


Fig. 17

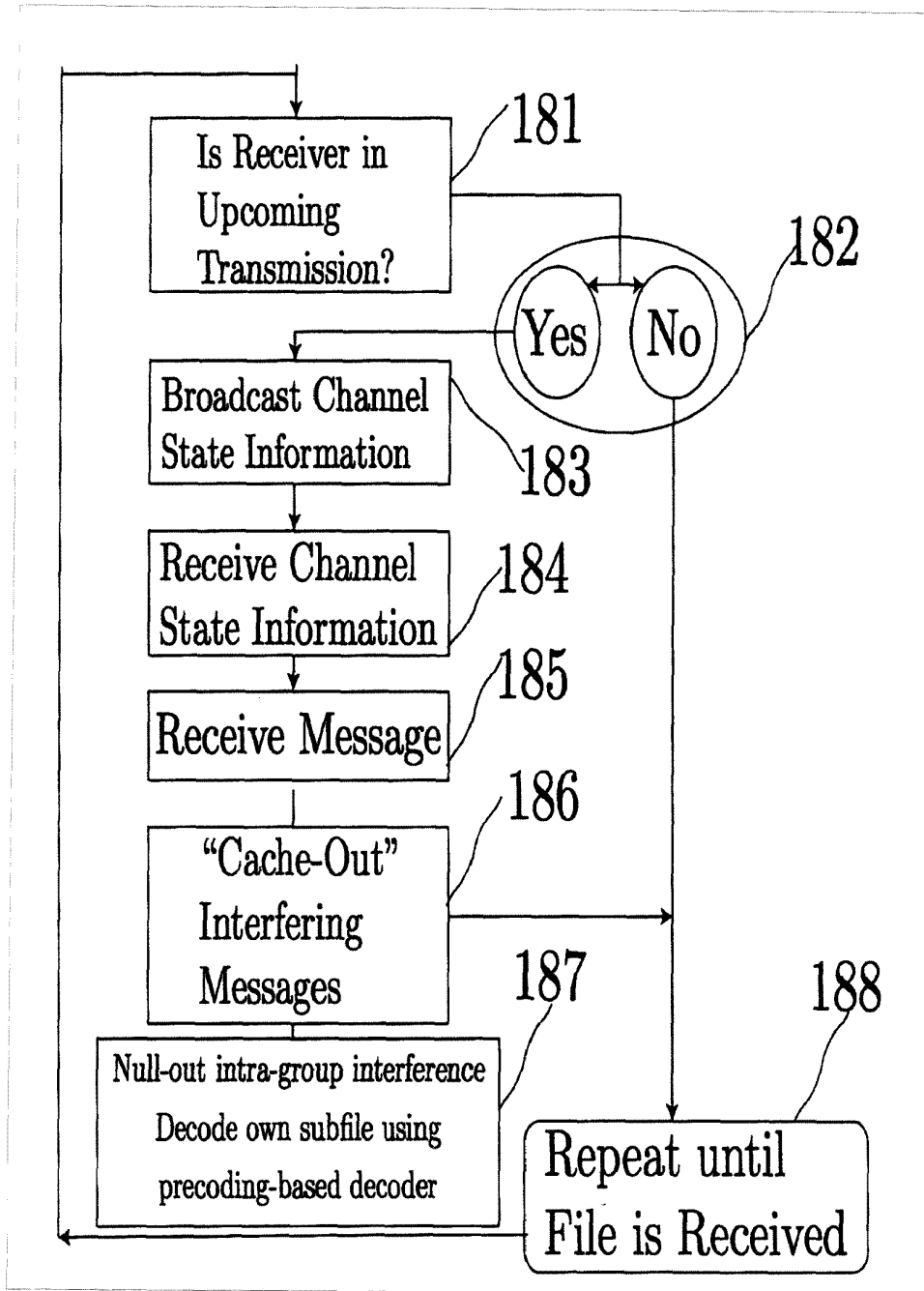


Fig. 18



EUROPEAN SEARCH REPORT

Application Number
EP 17 29 0111

5

DOCUMENTS CONSIDERED TO BE RELEVANT			
Category	Citation of document with indication, where appropriate, of relevant passages	Relevant to claim	CLASSIFICATION OF THE APPLICATION (IPC)
X	SEYED POOYA SHARIATPANAH ET AL: "Multi-Antenna Coded Caching", ARXIV.ORG, CORNELL UNIVERSITY LIBRARY, 201 OLIN LIBRARY CORNELL UNIVERSITY ITHACA, NY 14853, 11 January 2017 (2017-01-11), XP080740868, * the whole document *	1,3-12, 16	INV. H04L29/08 H04B7/0413 H04B7/024
Y A	-----	2 13-15	
X	SHARIATPANAH SEYED POOYA ET AL: "Multi-Server Coded Caching", IEEE TRANSACTIONS ON INFORMATION THEORY, IEEE PRESS, USA, vol. 62, no. 12, 1 December 2016 (2016-12-01), pages 7253-7271, XP011634636, ISSN: 0018-9448, DOI: 10.1109/TIT.2016.2614722 [retrieved on 2016-11-18]	13-15	
A	* page 1 - page 9 * * page 12 - page 18 * * sections I, II, III, V * * appendices B, C * * figures 1,2 * ----- -/--	1-12,16	TECHNICAL FIELDS SEARCHED (IPC) H04L H04B
The present search report has been drawn up for all claims			
Place of search Munich		Date of completion of the search 29 May 2018	Examiner Kokkinos, Titos
CATEGORY OF CITED DOCUMENTS X : particularly relevant if taken alone Y : particularly relevant if combined with another document of the same category A : technological background O : non-written disclosure P : intermediate document		T : theory or principle underlying the invention E : earlier patent document, but published on, or after the filing date D : document cited in the application L : document cited for other reasons & : member of the same patent family, corresponding document	

EPO FORM 1503 03.82 (P04C01)

10

15

20

25

30

35

40

45

50

55



EUROPEAN SEARCH REPORT

Application Number
EP 17 29 0111

5

DOCUMENTS CONSIDERED TO BE RELEVANT			
Category	Citation of document with indication, where appropriate, of relevant passages	Relevant to claim	CLASSIFICATION OF THE APPLICATION (IPC)
X	LIU AN ET AL: "Cache-Enabled Opportunistic Cooperative MIMO for Video Streaming in Wireless Systems", IEEE TRANSACTIONS ON SIGNAL PROCESSING, IEEE SERVICE CENTER, NEW YORK, NY, US, vol. 62, no. 2, 1 January 2014 (2014-01-01), pages 390-402, XP011536443, ISSN: 1053-587X, DOI: 10.1109/TSP.2013.2291211 [retrieved on 2013-12-31]	16	
A	* page 1 - page 4 * * section I - section II; figures 1-3 *	1-15	
Y	----- APOSTOLOS DESTOUNIS ET AL: "Alpha Fair Coded Caching", ARXIV.ORG, CORNELL UNIVERSITY LIBRARY, 201 OLIN LIBRARY CORNELL UNIVERSITY ITHACA, NY 14853, 26 January 2017 (2017-01-26), XP080751800, * page 1 - page 4 * * section I - section III; figure 2 *	2	
A	----- MADDAH-ALI MOHAMMAD ALI ET AL: "Fundamental Limits of Caching", IEEE TRANSACTIONS ON INFORMATION THEORY, IEEE PRESS, USA, vol. 60, no. 5, 1 May 2014 (2014-05-01), pages 2856-2867, XP011545810, ISSN: 0018-9448, DOI: 10.1109/TIT.2014.2306938 [retrieved on 2014-04-17] * the whole document *	1,3-16	
A,D	----- MADDAH-ALI MOHAMMAD ALI ET AL: "Fundamental Limits of Caching", IEEE TRANSACTIONS ON INFORMATION THEORY, IEEE PRESS, USA, vol. 60, no. 5, 1 May 2014 (2014-05-01), pages 2856-2867, XP011545810, ISSN: 0018-9448, DOI: 10.1109/TIT.2014.2306938 [retrieved on 2014-04-17] * the whole document *	1-16	TECHNICAL FIELDS SEARCHED (IPC)
The present search report has been drawn up for all claims			
Place of search Munich		Date of completion of the search 29 May 2018	Examiner Kokkinos, Titos
CATEGORY OF CITED DOCUMENTS X : particularly relevant if taken alone Y : particularly relevant if combined with another document of the same category A : technological background O : non-written disclosure P : intermediate document		T : theory or principle underlying the invention E : earlier patent document, but published on, or after the filing date D : document cited in the application L : document cited for other reasons ----- & : member of the same patent family, corresponding document	

EPO FORM 1503 03.82 (P04C01)

10

15

20

25

30

35

40

45

50

55



Application Number

EP 17 29 0111

5

CLAIMS INCURRING FEES

The present European patent application comprised at the time of filing claims for which payment was due.

Only part of the claims have been paid within the prescribed time limit. The present European search report has been drawn up for those claims for which no payment was due and for those claims for which claims fees have been paid, namely claim(s):

No claims fees have been paid within the prescribed time limit. The present European search report has been drawn up for those claims for which no payment was due.

10

15

20

LACK OF UNITY OF INVENTION

The Search Division considers that the present European patent application does not comply with the requirements of unity of invention and relates to several inventions or groups of inventions, namely:

see sheet B

All further search fees have been paid within the fixed time limit. The present European search report has been drawn up for all claims.

As all searchable claims could be searched without effort justifying an additional fee, the Search Division did not invite payment of any additional fee.

Only part of the further search fees have been paid within the fixed time limit. The present European search report has been drawn up for those parts of the European patent application which relate to the inventions in respect of which search fees have been paid, namely claims:

None of the further search fees have been paid within the fixed time limit. The present European search report has been drawn up for those parts of the European patent application which relate to the invention first mentioned in the claims, namely claims:

The present supplementary European search report has been drawn up for those parts of the European patent application which relate to the invention first mentioned in the claims (Rule 164 (1) EPC).

25

30

35

40

45

50

55



LACK OF UNITY OF INVENTION
SHEET B

Application Number
EP 17 29 0111

5

The Search Division considers that the present European patent application does not comply with the requirements of unity of invention and relates to several inventions or groups of inventions, namely:

10

1. claims: 1-12, 16

A process of communication for the system of D1, where caching is at least a function of said grouping.

15

2. claims: 13-15

A process of communication between K independent nodes N₁, N₂, ..., N_K, where said nodes are requesting at least one out of a plurality of files from a known library and where at least some nodes have caching capabilities and wherein a number of the transmitting nodes are cooperating for simultaneous transmission wherein the process determines caching at the nodes as a function of, among other things, the maximum number of file segments per file.

25

30

35

40

45

50

55

REFERENCES CITED IN THE DESCRIPTION

This list of references cited by the applicant is for the reader's convenience only. It does not form part of the European patent document. Even though great care has been taken in compiling the references, errors or omissions cannot be excluded and the EPO disclaims all liability in this regard.

Non-patent literature cited in the description

- **L. ZHENG ; D. N. C. TSE.** Communication on the Grassmann manifold: A geometric approach to the noncoherent multiple-antenna channel. *IEEE Trans. Information Theory*, February 2002 [0011]
- **M. MADDAH-ALI ; U. NIESEN.** Fundamental limits of caching. *IEEE Trans. Information Theory*, May 2014 [0011]
- **K. SHANMUGAM ; M. JI ; A. M.TULINO ; J. LLORCA ; A. G. DIMAKIS.** Finite length analysis of caching-aided coded multicasting. *ArXiv e-prints*, August 2015 [0011]
- **S. P. SHARIATPANAHI ; A. S. MOTAHARI ; B. H. KHALAJ.** Multi-server coded caching. *ArXiv e-prints*, August 2015 [0011]
- **N. NADERIALIZADEH ; M. A. MADDAH-ALI ; A. S. AVESTIMEHR.** Fundamental limits of cache-aided interference management. *CoRR*, vol. *abs/1602.04207*, 2016, <http://arxiv.org/abs/1602.04207> [0011]