



HAL
open science

Vehicular Carriers for Big Data Transfers

Raul Adrian Gorcitz, Yesid Jarma, Prométhée Spathis, Marcelo Dias de Amorim, Ryuji Wakikawa, John Whitbeck, Vania Conan, Serge Fdida

► **To cite this version:**

Raul Adrian Gorcitz, Yesid Jarma, Prométhée Spathis, Marcelo Dias de Amorim, Ryuji Wakikawa, et al.. Vehicular Carriers for Big Data Transfers. 2012 IEEE Vehicular Networking Conference (VNC), Nov 2012, Seoul, South Korea. pp.1-8. hal-00739361v1

HAL Id: hal-00739361

<https://hal.science/hal-00739361v1>

Submitted on 8 Oct 2012 (v1), last revised 9 Oct 2012 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Vehicular Carriers for Big Data Transfers

Raul A. Gorcitz*, Yesid Jarma*, Prom  th  e Spathis*, Marcelo Dias de Amorim*,
Ryuji Wakikawa , John Whitbeck*, Vania Conan , and Serge Fdida*

*UPMC  Toyota ITC  Thales

Abstract—In the latest years, Internet traffic has increased at a significantly faster pace than its capacity, preventing efficient bulk data transfers such as datacenter migrations and high-definition user-generated multimedia data. In this paper, we propose to take advantage of the existing worldwide road infrastructure as an offloading channel to help the legacy Internet assuage its burden. One of the motivations behind our work is that a significant share of the Internet traffic is elastic and tolerates a certain delay before consumption. Our results suggest that piggybacking data on vehicles can easily lead to network capacity in the petabyte range. Furthermore, such a strategy exceeds by far the performance of today’s alternatives that, although yielding good performance levels, still rely on the legacy Internet and inherent then its intrinsic limitations. We show through a number of analyses that our proposal has the potential to obtain remarkable reductions in transfer delays while being economically affordable.

I. INTRODUCTION

During the last ten years, the estimated total number of Internet users on the planet increased from 500 millions to 2.4 billions [1]. Recent statistics show that, in a single day, 22 million hours of TV shows are watched on Netflix, 864,000 thousand videos are uploaded on Youtube, 18.7 million hours of music are streamed on Pandora, and the amount of information that transits over the Internet is enough to fill 168 million DVDs [2]. According to Cisco, global IP traffic has increased eight times in the last five years and is expected to increase four times in the next five years [3]. Commonly referred to as the exaflood and the information explosion, the rapidly increasing amount of traffic transferred over the Internet is largely driven by data-intensive applications.

Motivated by the need for technical flexibility and cost-effective scalability, large companies, organizations, universities, and governmental agencies constantly move their data and applications within and between data centers to balance workloads, handle replication, and consolidate resources. As a result, the demand for bandwidth-intensive services such as cloud computing, multimedia transfers, data migration, disaster recovery, and online backups has strained the Internet infrastructure to its limits.

Despite the ever-growing demand in bulk transfers, the price of bandwidth remains prohibitively high, especially at the network core. As a result, many edge providers are rate-limiting or even blocking the use of bandwidth-intensive applications. While bulk traffic can be seen as expensive when considering the high bandwidth consumption incurred, data-intensive applications are also less demanding when it comes to the requirements in terms of delay. Compared to most interactive applications that are highly delay-sensitive, the average

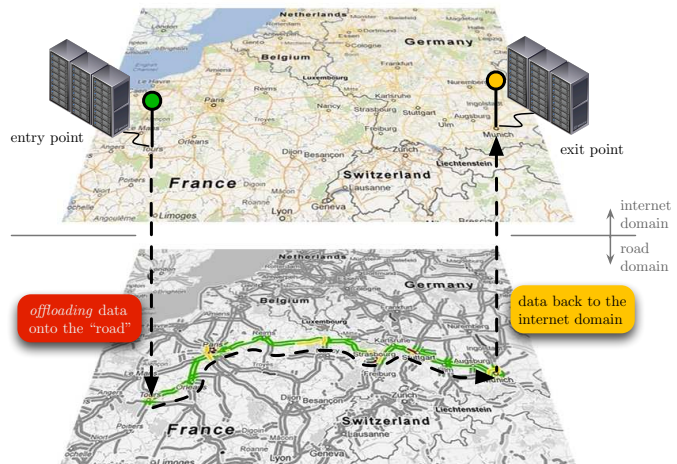


Fig. 1. At some point in space, data is offloaded from the Internet domain onto the road domain and is carried by vehicles up to some destination address, saving significant Internet resources.

throughput is the main criterion to evaluate the performance of bulk transfers and also to improve user experience.

Current methods for transferring such data include adaptations to standard file transfer methods or the use of hard drives and DVDs together with a courier service [4]. These solutions, although simple, can be either time consuming, or costly, or both. Recent alternatives propose to schedule data traffic to off-peak hours or using transit storage nodes [5], [6]. Nevertheless, as long as the legacy Internet is used, the underlying technology stays the same and the ISPs are still handling all the traffic. Cho and Gupta suggest to combine the Internet with the postal system to send a part of the data using hard-drives [7]. Nevertheless the bandwidth consumption requirements are on a steady rise and are mostly dictated by demands at peak hours.

We argue that disruptive solutions should be considered when it comes to moving huge amounts of data between geographically distributed sites. In this paper, *we propose to exploit the delay-tolerant nature of bulk transfers to deliver data over the existing road infrastructures*. Our work is motivated by the increasing number of vehicles driven and miles traveled in the world. The vehicle fleet in operation worldwide surpassed the 1 billion mark in 2010 and is expected to double in the next two decades. The number of vehicles ownership is forecast to grow to up to 4 billion by mid-century. By leveraging the communication and storage capabilities that will

soon equip (if not already) vehicles, *we advocate the use of conventional vehicles as the communication medium for big data migration in an opportunistic manner.*

The idea is illustrated in Fig. 1. Given the flow of vehicles daily traveling roads, our system design built on top of the road network can effectively *offload* the legacy Internet infrastructure for massive delay-tolerant data transfers. Furthermore, we suggest that a company with a well developed business model can provide the right incentives for vehicle owners to become a part of such a system. To evaluate our system, we compare our results with a state-of-the-art bulk data transfer scheme and we show that a vehicular-based solution may lead to significant improvements in terms of bandwidth while achieving the right tradeoff between transfer delay and cost.

In summary, the contributions of our work are as follows:

- **Offloading scheme.** We describe a novel vehicular-based opportunistic bulk data transfer system, designed to offload the Internet from delay-tolerant content.
- **Capacity improvement.** We compute the transfer delay for vehicular-based bulk data transfers, and we show that the system has the potential of moving massive amounts of data in short time periods compared to today's traditional techniques.
- **Reduced cost.** We evaluate the cost of such a system by suggesting a business model meant to motivate drivers to become a part of it.

The remainder of this paper is structured as follows. In Section II, we describe the overall system operation and list the assumptions we consider in this work. In Sections III and IV, we evaluate the potential of our solution in terms of transfer delay and cost. We give insights into some interesting open issues in Section V and postpone the related work to Section VI so that the reader has enough material to better understand the positioning of our work with regard to the literature. We finally conclude the paper in Section VII.

II. OFFLOADING ONTO CONVENTIONAL VEHICLES

In this body of work we argue that by taking advantage of the characteristics of future *smart* vehicles, such as data storage capabilities, these can be used to transport massive quantities of data between two geographical locations. By using the highly developed worldwide road and highway infrastructure, vehicles can offload high volumes of data from the Internet.

A. System operation

In order to overcome the limitations in terms of capacity and design of the Internet, vehicles are equipped with one or more removable memory storage devices such as magnetic disks or other non-volatile solid-state storage devices. The term “vehicle” refer to both passenger and commercial vehicles; in the latter case, it may be part of a fleet vehicle owned or leased by a business or governmental agency. We assume that vehicles also embed one or more communication network interfaces and a positioning system. The system we describe

below includes vehicles in operation and their users, a service provider, and a content provider.

Memory devices can be owned by a party other than the user of the vehicle. Typically, a service provider may own the memory devices and owners of the vehicles can be compensated based on the amount of data transferred. A content provider distributes the data to be piggybacked onto the vehicles through a wide-area data network such as the Internet. The service provider charges the content provider for the amount of data to be transferred along the road infrastructure. A network of *offloading spots* provides the data to be piggybacked on the memory of vehicles. We use the term “offloading spots” to refer to locations that provide the data to be transferred to the memory devices of the vehicles or where on-board memory devices can be exchanged for pre-loaded memory devices that match the destination of the vehicle. The offloading spots can be placed at locations where vehicles may be parked. For example, an offloading spot can be located in a shopping center parking lot, a street parking spot, or at the users' home place. At a higher level, a collection of offloading spots form an entry/exit point as illustrated in Fig. 1.

The service provider selects the offloading spots from the group consisting of the loading stations that transfer the data to the already in-place memory of the vehicle and the memory swap stations that replace the on-board memory of the vehicle. Vehicle memories can be loaded with data while parking at the offloading spot or exchanged for ready-to-ship memory devices so that users can continue their travels without waiting for the data to be loaded. The selection of the offloading spots is based in part on the geographic location of the vehicle and if available, its planned destination. The service provider also monitors the status of the offloading spots which include the available parking space, the memory exchange bays that are free, the destination of the data made already available for shipment. The service provider also periodically queries the vehicles over the data network to determine the current geographic location and destination of the vehicles. The positioning system of the vehicle includes a navigation system that generates routes and guidance between a geographic location and a destination. The historical locations and addresses are stored in a geographic location database managed at the service provider's control center. The service provider also keeps record of the status of the offloading spots in a specific database. The service provider matches the destination of the vehicles to a group of offloading spots selected based on park space availability. If preloaded memory devices are ready to be shipped, the service provider checks for free exchange bays at the offloading spots or contacts the content provider in order to transfer the data to be loaded on the memory devices.

B. Assumptions

Let us first denote the frequency of vehicles passing an entry point and traveling towards the same exit point (destination) as f . This value will be important to compute the maximum achievable capacity of the system. In practice, we assume that not all vehicles will have enough incentives to take part of the

TABLE I
SUMMARY OF THE VARIABLES USED THROUGHOUT THIS PAPER.

symbol	meaning
τ	transmission delay for the entire data
\mathcal{D}	total data to be transferred
f	vehicle frequency at a point
\mathcal{S}	storage capacity of a vehicle
\mathcal{P}	penetration ratio of the technology
d	travel distance
\bar{s}	average speed
s_f	free-flow speed
K	traffic jam vehicle density
k	actual vehicle density

system. We call \mathcal{P} the penetration ration of the technology, i.e., the probability that a vehicle accepts to carry data. We set this value to 20%, which corresponds to the approximate value of the market share of Toyota in California [8].¹ We assume that all the vehicles that enter the highway at entry point A reach exit point B . Therefore, the value of $f \times \mathcal{P}$ for vehicles leaving A is identical for vehicles reaching B . This allows us to calculate the vehicular system performance unhindered by routing errors and data loss. Although we are aware that in a real life scenario a certain number of vehicles might exit the highway before reaching point B , the purpose of this work is to present and demonstrate the potential of performing vehicular-based bulk data transfers in an opportunistic way and therefore we do not discuss data loss and routing related issues. We invite the reader to refer to Section V for some discussion regarding this subject.

We denote the total amount of data to be transferred between A and B as \mathcal{D} . This amount is chunked and divided among the participating vehicles. For this computation, we denote as \mathcal{S} the storage capacity of each vehicle. All this gives $f \times \mathcal{P} \times \mathcal{S} = \mathcal{D}$. Finally, we call d the travel distance between A and B and \bar{s} the average speed of the highway. The summary of the variables used in this paper is shown in Table I.

Note that the proposed approach would be worth deploying only if its performance exceeds the one obtained by other means or if its value-add is confirmed. Therefore, in order to better evaluate the benefits and possible pitfalls of the proposed system we use the following performance metrics: *transfer latency*, *system throughput*, and *financial cost*. The first two will be detailed in Section III while the latter will be discussed and evaluated in Section IV.

III. VEHICULAR CARRIERS: PERFORMANCE ANALYSIS

We evaluate the performance of bulk data transfers using vehicular carriers as described in the previous section.

A. Dataset

Our results are based on a vehicular frequency dataset spanning a two year period, made publicly available at the

¹Note that this choice heavily impacts our results. However, we believe that this would be a lower bound, which corresponds to the case where a single car manufacturer accepts to embed this technology on their vehicles. Real values would be hopefully better than the numbers we show in this paper.

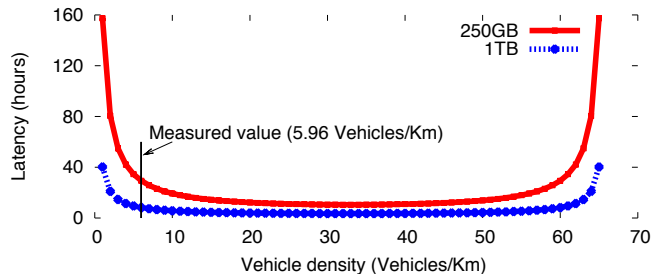


Fig. 2. Total delay required to transfer 1 Petabyte as a function of the highway density. The vertical line represents the measured value obtained from the dataset for the highway between Orleans and Tours (see more results below).

end of 2011 by the French Ministry of Ecology, Sustainable Development, and Energy [9]. In this dataset, vehicular frequency measurements were conducted on multiple segments of the highway, yielding slightly different results at each measurement (as vehicles enter and exit the highway segment). Since we consider that all vehicles starting a journey from point A will reach point B , to filter out vehicles exiting the road in between, we use the minimum vehicle flow value measured on the trajectory. For the sake of clarity, we focus our analysis on a stretch of highway between two adjacent cities in France, namely Orleans and Tours. The results for other scenarios are presented later in this section. The stretch of highway under consideration is 118-Km long and has 3+1 lanes.

The data is provided under the AADT (Annual Average Daily Traffic) engineering standard used for vehicle traffic load on a given section of road. This type of standard is used for transportation planing/engineering and traffic related pollution [10], [11]. The AADT is an average measured for a period of one year and divided by the number of days. This type of data is useful in avoiding traffic differences depending on season or time of the day.

B. Computing transfer delays

The time required to transport data along the highway between two locations can be calculated as a sum of two factors. The first one is the time required for a vehicle to transit the stretch of highway between the starting point and its destination. The second one is the amount of time required to load all the information onto vehicles. The total transfer time depends on the storage capacity of each vehicle, on the penetration ratio of the vehicular bulk data transfer technology, and on the frequency of vehicles.

The vehicle frequency measurement is one of the core elements required for capacity calculations in the highway research area [12], [13]. Highways are designed to avoid traffic jams and facilitate the free flow of traffic by estimating the vehicle frequency in a stretch of highway and the maximum vehicular density (i.e., the number of vehicles per unit of distance). Since for our proposal traffic jams could severely hinder the performance of the system, our calculations are based on the vehicle density. The vehicular frequency f at

one point on the highway, can be expressed in terms of the vehicular density as follows:

$$f = s_f \times \left(k - \frac{k^2}{K} \right), \quad (1)$$

where s_f is the free-flow speed (i.e., maximum speed), K is the density of vehicles characterizing a traffic jam, and k is the actual density measured in vehicles per kilometer (a traffic jam takes place when the value of k approaches K). Therefore, the total transfer latency τ can be expressed as a function of the vehicular density as follows:

$$\tau = \frac{D \times K}{s_f \times ((k \times K) - k^2) \times \mathcal{S} \times \mathcal{P}} + \frac{d}{\bar{s}}. \quad (2)$$

In Fig. 2, we present the transfer latency for transporting 1 PB of data in function of the density of vehicles using Equation 2. We consider vehicles with storage capacities of 250 and 1,000 Gigabytes, with a class A level of service on a rural highway with 3 lanes [12].² On one extreme, long delays are obtained with low traffic densities as less vehicles are available to act as carriers. As density increases, we observe a significant increase in performance until we reach the optimal density. On the other extreme, as vehicle density grows beyond the optimal value, congestion levels rise and the speed of the vehicle flow slowly decreases until it reaches a jam zone that causes a steep increase of the latency parameter.

In Fig. 2, the vertical line indicates the density obtained from the dataset (5.96 vehicles/Km), which is the reference value we will use throughout the rest of the paper. Note that, in this case, even though the measured flow is largely below the optimal capacity of the highway, the transfer latency values obtained are as low as 8 hours to transfer 1 PB of data, and has the potential of being lower as the vehicle density approaches the design optimum.

In Fig. 3(a), we show the transfer latency as a function of the total amount of data to be transferred. The two curves represent vehicles equipped with 250-GB and 1-TB storage unities. Our proposal is able to obtain delays of under 9 hours for quantities of data running up to 1 PB using a storage capacity of 1 TB. As we will see later in this paper, these values overcome by far the performance obtained by alternative solutions that rely on the current Internet architecture. In Fig. 3(b), we show the throughput of the system in function of the total data to be transferred. We also derive the theoretical throughput of the system and show the results in Fig. 3(b).

Other highways. We have performed the same calculations for a number of highways in France and show the results Fig. 4. As we can see, the values shown in the figure are compliant with the numbers obtained for Orleans \leftrightarrow Tours.

²The ‘‘class A’’ level of service refers to very good driving conditions where the drivers are unhindered by other traffic participants and are able to maintain the desired speed. Vehicle density values are the main performance metric used for the level of service estimation.

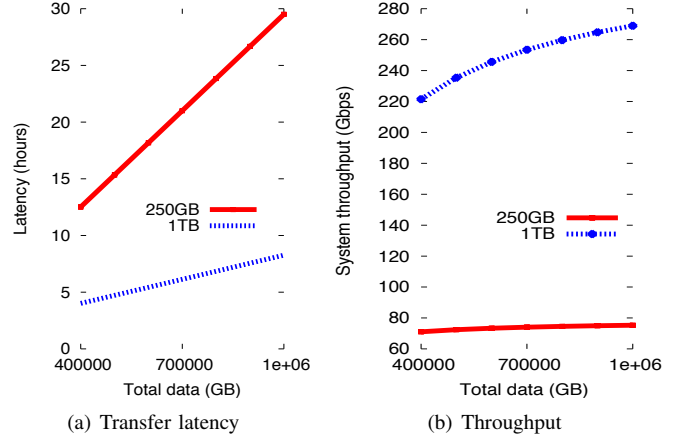


Fig. 3. Transfer latency and throughput of the system for the road connecting Orleans and Tours. We consider different values of D and S . Note that the curves are not linear because of the variable term in Equation 2.

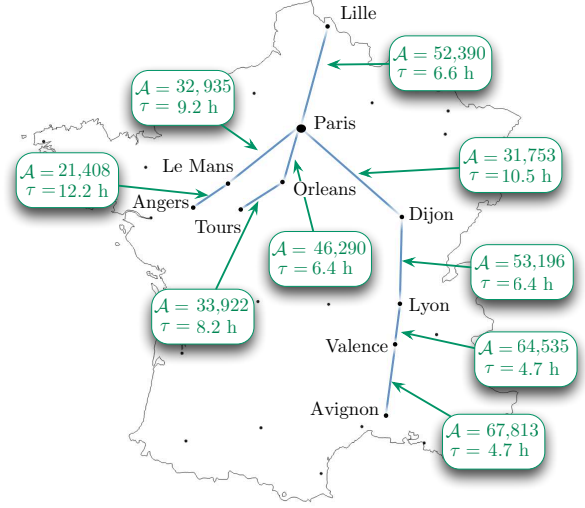


Fig. 4. Average Transfer Delays (τ) obtained using the ‘‘annual average daily traffic’’ (\mathcal{A}) on segments of highways connecting several important locations in France. The parameters used here are: average speed = 100 Km/h, total data = 1 PB, per-vehicle storage capacity = 1 TB, and penetration ratio = 20%.

C. Impact of ‘‘environmental parameters’’

As we have seen from Equation 2, other than the vehicle density and the amount of data to be transferred, the performance of our system depends as well on other ‘‘environmental’’ factors such as distance and speed. Therefore, it is necessary to assess the impact of both these factors.

1) *Average speed \bar{s} :* In order to better understand the impact of the average vehicle speed, we have considered the same parameters as before with the exception of the total data to be transferred, which has been fixed at 1 PB, and the data storage capacity which has also been fixed at 1 TB. We vary the average speed from 60 Km/h to 130 Km/h. The results are shown in Fig. 5. It is clear that, even if the average speed has a non-negligible impact on the performance of the solution, the system is relatively impervious to changes in the average

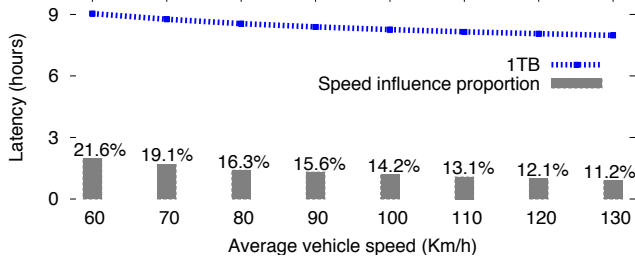


Fig. 5. System performance in terms of latency, affected by the variability of the speed parameter and the proportion of its impact. The total data to be transferred is set to 1 PB with a vehicle storage capacity of 1 TB.

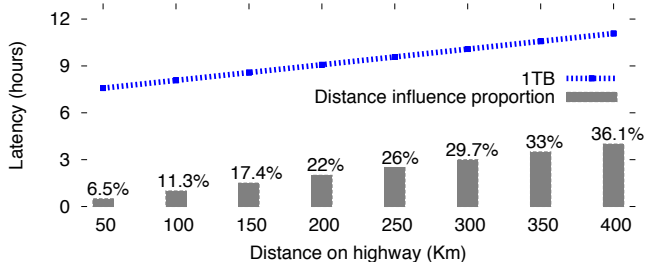


Fig. 6. System performance in terms of latency, affected by the variability of the distance parameter and the proportion of its impact. Total data to be transferred is set to 1 PB with a vehicle data storage capacity of 1 TB.

speed on the highway.

2) *Distance d*: Because the length of highways linking two points of interest can differ significantly, we have tested our proposed solution on longer distances by varying them from 50 Km to 400 Km. As in our previous experiment, we have fixed the per-vehicle storage capacity to 1 TB, the average speed to 100 Km/h, and the total transferable data to 1 PB, while keeping all the other parameters unchanged. The results are depicted in Fig. 6.

D. Comparison with NetStitcher

In this section we compare the vehicular-based bulk data transfer system with an Internet-based bulk data transfer solution named NetStitcher [6] (see also Section VI). NetStitcher achieves its best performance by transferring 1.15 Terabytes in 3 hours between cities in the same time-zone. These values will be used as a hard constraint during our comparison scenario. Here we are asking the following question: *is our vehicular-based bulk data transfer system able to transfer a comparable amount of data in the same period?* To answer this question, we developed a scenario that uses the segment of highway between two cities in the same time-zone, namely Orleans and Tours. We have varied the data storage capacity of the vehicles between 250 GB and 1 TB, obtaining two distinct performance values for the amount of data transferable in a 3-hour time interval. The results are presented in Fig. 7. The vehicular-based system outperforms by far the Internet-based bulk data transfer system by transferring up to 200 times more data in the same amount of time.

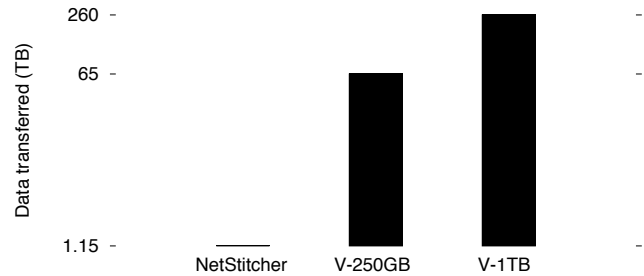


Fig. 7. Comparison between an Internet-based solution (NetStitcher [6]) and the proposed vehicular-based solution (V- 250GB and V-1TB), depicting the amount of data transferable in a delay of 3 hours between two cities in the same timezone.

IV. FINANCIAL COST EVALUATION

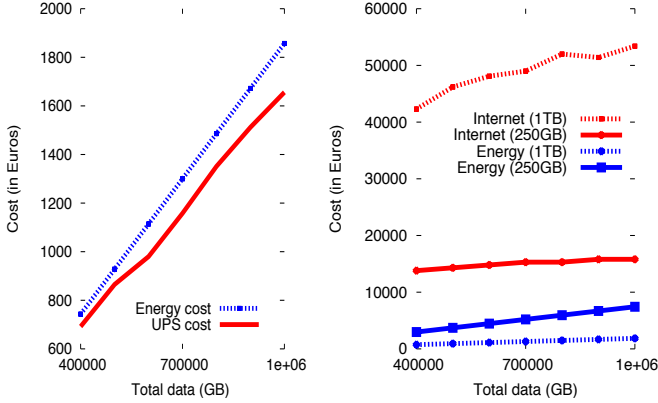
In this section, we discuss the financial cost of transporting massive amounts of data in an opportunistic manner using vehicles as data carriers. For comparison purposes, we have imagined a business model in order to motivate vehicle owners to participate in the system. We compare the cost of our model with both a package delivery solution and an Internet-based solution.

A. Comparison with a package delivery service

We consider that a company or a group of associated companies need to migrate 1 PB of data from two different data centers located in Orleans and Tours respectively. To this end, companies use a package delivery service such as UPS [14]. The total cost in such a case is composed of mainly two parts, the acquisition of transportable storage equipment and the cost of transporting the equipment to the destination. For this purpose, we have chosen hard drives with the highest storage capacity available in the market.³ Seagate Constellation ES 2 has a 3 TB capacity and a mass of 700 grams [15]. The total amount of data to be sent was divided between the capacity of the hard drives and packed into boxes that could contain up to 40 hard-drives. Knowing the weight and the size of the boxes, we were able to calculate the delivery price using the company's own on-line application. The packages have an estimated delivery time between 19 and 91 hours, with the fastest service available. The difference in delay is given by the company's own shipping schedule (e.g., a package sent on Friday will only be delivered on Tuesday the week after).

Next, we have considered that a company decides to use the vehicular-based data transfer system. It conceived a hypothetical business model in which the company pays for the electricity cost of recharging the vehicle's batteries at a recharge station, in order to motivate drivers to participate. As it is hard to predict the average range and consumption

³The data is transported over a large area and fragile storage equipment, such as magnetic tapes are exposed to environmental factors like humidity, heat or magnetic interference which could potentially damage the magnetic tape. Other storage devices such as Blu-ray Discs have also been excluded as they have smaller storage capacities and data transfer rates. In order to obtain the best compromise between storage capacity, performance and robustness, we have chosen hard drives as our transportable data storage equipment.



(a) Cost comparison between the electrical vehicles full recharge cost for traveling from Orleans to Tours and the transport cost using a package delivery service between the same destinations. (b) Cost comparison between an Internet-based implementation using dedicated links and a vehicular-based implementation using vehicle carriers. Vehicle storage capacity has been alternated between 250GB and 1TB.

Fig. 8. Cost comparison between UPS and our proposed approach.

of electric vehicles in the next 20 years, we decided to use the technology available today. For this purpose, we have considered an electric car with a battery that has a capacity of about 22 kWh and range of 180 Km. In accordance to this study on Electric Vehicles [16], 75% of electric vehicles in 2011 had a maximum autonomy of 200 Km and an average battery capacity of 23.4 kWh. This gives us an average consumption of about 0.12 kWh per kilometer traveled. We have also considered the price of 0.1312 € per kWh, as being paid by a non-commercial customer in France during peak hours [17]. We first compute the number of vehicles required to transport the data, which is $N = \mathcal{D}/S$. We can then easily calculate the price of a single full recharge:

$$c_{\text{unit}} = d \times U \times E, \quad (3)$$

where d is the distance travelled, U is the average electricity consumption of a vehicle, and E is the price of electricity per kWh. The total energy cost is then:

$$c_{\text{total}} = N \times c_{\text{unit}}. \quad (4)$$

We then compare the costs of the package delivery and the vehicular-based system. The results are presented in Fig. 8(a). The plot shows the cost of the total recharge price for a given amount of total data to be transported considering 1 TB of per-vehicle storage capacity. The results show similar performance in terms of cost yet a huge difference exists in terms of delay as transporting 1 PB of data with vehicles can take less than 4 hours while using UPS takes a minimum of 19 hours. Considering the fact that certain types of delay tolerant content have very strict constraints in terms of latency, an addition of 15 hours to the transfer time could render the data unusable.

B. Comparison with an Internet-based solution

We now discuss the cost of massive data transfers that use the current Internet infrastructure by means of dedicated links.

The latency results obtained in Section III are a major performance indicator and will be used as data transfer constraints in this section.

In our calculations, the theoretical throughput achieved by the vehicular-based system is not influenced by the inefficiencies of an Internet-based technology. Thus the values presented in the previous section are the theoretical throughput of the vehicular system at 100% efficiency. We are well aware that in a TCP Internet-based transfer system, the actual data transfer capacity is limited. We do not debate on the way the dedicated links should be managed in order to obtain the equivalent TCP throughput of the vehicular system. Here we assume the entire bandwidth is used efficiently.

In Fig. 8(b), we depict two scenarios represented by two different curves each. The two curves at the top represent the cost of purchasing dedicated links so as to match the vehicular-based system performance in terms of latency and throughput. In order to correctly compute the values,⁴ we have used the monthly dedicated links price offers from a major Point of Presence in France [18]. The two curves at the bottom represent the cost of electricity required to recharge the vehicle carriers after delivering their load. Here the cost is also calculated according to the same latency and throughput constraints as in the previous scenario.

We can easily evidence a major difference in terms of cost. The poor performance of the Internet-based system is due to a latency requirement that is hard to obtain for massive data transfers over the Internet. A large number of dedicated links are being used to transfer the given amount of data to be able to respect these latency constraints. The efficiency of such an implementation could be debated as the price paid for the dedicated links is hard to justify when massive data transfers are rarely made during a month. When requiring shorter delays to transfer massive amounts of data, the entire capacity of the dedicated links is only used during short periods of time, while during the rest of the time resources are being used inefficiently.

V. DISCUSSION

The goal of this section is to point out several research topics that could contribute to the development of this research topic.

Routing and data delivery. One of the most pressing issues that requires attention is data delivery over longer distances. In the case of a highway linking three consecutive cities, it is not reasonable to make the assumption that all vehicles leaving the first city will reach the third city. This motivates the development of specialized routing protocols that limit data loss and ensure efficiency when drivers stop at intermediary points. The limited battery capacity would force the driver to stop for example at recharge stations or at *battery swap stations*. Such “points of interest” could act as routers where data could be swapped from one vehicle to another, heading in the right direction.

⁴The price was calculated by aggregating the highest available links.

Scheduling and transfer planning. Vehicular traffic has observable diurnal patterns with dramatic increases in vehicular frequency during rush hours contrasting with much smaller values during night time when most of the traffic is represented by commercial freighters. This pattern of movement is dictated by constraints that characterize the driver’s behavior (like working hours or driving preference) and it can be anticipated to a certain extent. In order to use efficiently the opportunistic vehicular-based system, data transfers must be properly scheduled and thus appropriate methods should be developed.

Data loading. Today’s technology is advancing and is offering new alternatives to loading data on to vehicles. If battery swap stations are deployed at a large scale, vehicles could be “fed” with data when they get a newly charged battery, and this in a very short period of time (about one minute [19]). Another option would be to use state-of-the-art microchips that allow wireless transmission speeds up to 1,000 times faster than current technology [20].

Incentives and business plans. Vehicles participating in the system would be private property and as such, they cannot be used to transfer data without the consent of the owner. A mutual beneficial business plan needs to be developed to motivate car owners to participate, while a partnership between multiple companies could potentially reduce the servicing costs of such a system.

VI. RELATED WORK

In this section, we first present some relevant data-intensive applications that motivate the need for bulk transfer systems over or in replacement of the legacy Internet. We then discuss the related work from two aspects according to which we relate our work: one is the improvements proposed to the Internet so as to accommodate the transfers of bulk data. The other is the new opportunities emerging from recent and future advances within the framework of delay-tolerant networking.

A. Bulk data transfers over the Internet

Exchange of large scientific datasets, email archival and personal documents backup, and distribution of high-resolution movies are all examples of bulk data transfers over the Internet.

Scientific instruments such as the Large Hadron Collider (LHC) at CERN and the Laser Interferometer Gravitational Wave Observatory (LIGO) can generate tens of terabytes or even petabytes of data that need to be disseminated to remote collaborators or to computational capable centers. Some research projects have developed their own data management systems which include components to manage the data transfers over wide-area networks but also to schedule, monitor, and manage the placement of data and the execution of analysis jobs according to specific policies. The various protocols and tools developed to handle bulk transfers of scientific data over the Internet extend the standard FTP protocol so as to meet requirements such as fault tolerance or transfer concurrency. The most notable example is GridFTP which can achieve data transfers of up to a few TeraBytes per day.

Today’s large scale data centers are facing these same issues. With the cloud technology gaining ground, data centers are increasing in size demanding larger aggregate bandwidth requirements. To improve the services offered to users, data needs to be moved closer to the consumer along the provision based infrastructure. It can also require data restoration functions in a disaster scenario, keeping services running while problems are being solved or a distribution scheme that requires multiple geographic locations to work efficiently. In [21] the authors find that provisioning high levels of bandwidth with today’s existing techniques has a determinant impact on the development and maintenance budget of a company.

Though bulk data transfers are at the basis of some popular services such as high-definition multimedia content delivery, many ISPs are using scheduling, traffic shaping, and queue management techniques to limit the rate of bandwidth-intensive applications. An example of service affected by the ISPs policies is Netflix, who is responsible for about 29% of North American fixed internet access bandwidth utilization [22]. Initially launched as a DVD rent-by-mail company, Netflix began streaming movies online in 2007 in an attempt to avoid the postal costs for delivering DVDs by mail. To address the ISPs rate-limiting policies and the costs of serving content, recommendation algorithms are expected to be used in combination with peer-to-peer networking by big data service providers.

The idea behind the use of these algorithms is the creation of geographically logical clusters of subscribers defined by common interest. Thus, instead of pointing content requests to the providers’ central server, the P2P network of subscribers’ boxes will allow content to be served from other boxes in the same cluster.

To address the shortcomings faced by bulk data applications, Laoutaris et al. proposed NetStitcher, a proposal that exploits diurnal patterns of network traffic to schedule bulk transfers at times of low link utilization depending on the time zones [6]. Although this technique allows a good return on investments for dedicated lines, the physical medium is still limited due to the high cost and limited capacity of today Internet’s infrastructure. Furthermore, efficient scheduling decisions require real-time information regarding the network load and a transport layer needs to be designed so as to be able to send all delay tolerant traffic when spare bandwidth is made available.

B. Assisted DTN

Some of the issues we address in this work are also related to the Delay Tolerant Network (DTN) paradigm and more specifically to one of its variants, the Assisted DTNs (A-DTN). A-DTN architectures involve various data carriers or forwarders ranging from buses to airplanes to compensate for the lack of continuous connectivity by bridging otherwise disconnected nodes. By increasing connectivity opportunities, these special-purpose nodes referred to as data mules, message ferries, and throw-boxes intend to increase the available network capacity.

In [23], authors propose the use of battery-powered devices with on-board storage, processing, and radio capabilities called throw-boxes, which once placed in strategic points within the network can increase the capacity of the network. In contrast to throw-boxes which are fixed nodes, other approaches exploit the movements of nodes they can proactively control to enhance the throughput of the system. The difference between these approaches depends on the level of randomness introduced to the nodes mobility patterns. In [24], humans and animals moving in a sparse sensor network along random paths are used as data MULEs (Mobile Ubiquitous LAN Extensions) to gather sensing data opportunistically which are later delivered to some collection point. Rather than adopting total randomness, [25] exploits the non-random mobility patterns of special nodes called message ferries in order to collect messages and deliver them to their destination. Ferry nodes can either adjust their trajectory so as to meet up with the non-ferry nodes and exchange the data to be delivered or move according to specific predefined routes. In the latter case, non-ferry nodes with knowledge of the ferries routes need to move close to a ferry to communicate with it.

While the previously listed proposals target limited-sized areas, [26] propose to exploit airplane passengers boarding airline flights to bridge remote airports. Messages to be delivered are loaded onto the passengers' mobile devices at the airport depending on their destination while they are waiting for their flight. Unlike data MULEs, airplanes follow regular prescribed schedules and also differ from data ferries in that their routes cannot be controlled. Nevertheless flights destination can still be matched to the message destination. Simulation results show that under certain conditions, carrying data over scheduled flights can achieve a similar throughput as a single TCP connection as long as the amount of data to be transferred is equal to capacity of three DVDs.

Another work [7] suggested the combined use of the Internet together with the postal system to send a part of the data using hard-drives. Even though this system offloads some of the data onto a carrier other than the Internet, it still relies on a methodology that requires detailed scheduling, and can only be used for data able to tolerate delays of up to several days. Furthermore data on physical media such as tapes or removable disk drives have the downside of requiring manual handling at arrival.

VII. CONCLUSION

In this paper, we proposed a novel technique for big data transfers by offloading data from the Internet onto a network composed on conventional vehicles in an opportunistic fashion. Through extensive evaluation, we have shown that using vehicles as data carriers can be highly delay-efficient in today's context, where the amount of data to be transferred is increasing at a fast pace, and current technologies are becoming obsolete. Moreover, we have presented an incentive-based motivational approach that is meant to increase the penetration ratio of data carrying vehicles and discussed the costs of implementing such a system. Finally, we point to

several interesting research directions in this new area, hoping to motivate and persuade the research community to continue further research in this field.

REFERENCES

- [1] "International Telecommunication Union," <http://www.itu.int/ITU-D/ict/statistics/>.
- [2] "Internet World Stats," <http://www.internetworldstats.com/usage/use017.htm>.
- [3] "CISCO," http://www.cisco.com/en/US/solutions/collateral/ns341/ns525/ns537/ns705/ns827/VNI_Hyperconnectivity_WP.html.
- [4] W. Allcock, J. Bester, J. Bresnahan, S. Meder, P. Plaszczak, and S. Tuecke, "Gridftp: Protocol extensions to ftp for the grid," <http://www.ggf.org/documents/GFD.20.pdf>, 2003.
- [5] N. Laoutaris, G. Smaragdakis, P. Rodriguez, and R. Sundaram, "Delay tolerant bulk data transfers on the internet," in *ACM Sigmetrics*, Seattle, WA, USA, Jun. 2009.
- [6] N. Laoutaris, M. Sirivianos, X. Yang, and P. Rodriguez, "Inter-datacenter bulk transfers with netstitcher," in *ACM Sigcomm*, Toronto, Ontario, Canada, Aug. 2011.
- [7] B. Cho and I. Gupta, "Budget-constrained bulk data transfer via internet and shipping networks," in *International Conference on Autonomic Computing*, Karlsruhe, Germany, Jun. 2011.
- [8] "CNCDA California Auto Outlook Market Report, SQ 2012," <http://www.cncda.org/secure/GetFile.aspx?ID=2358>.
- [9] "Bison Futé," <http://www.bison-fute.equipement.gouv.fr/>.
- [10] *Traffic Monitoring Guide*. U.S. Department of Transportation, Federal Highway Administration, 2001.
- [11] T. Wright, P. S. Hu, J. Young, and A. Lu, "Variability in traffic monitoring data," <http://wwwwcf.fhwa.dot.gov/ohim/flawash.pdf>, U.S. Department of Transportation, Federal Highway Administration, Tech. Rep., Aug. 1997.
- [12] *Highway Capacity Manual 2010*. U.S. Transportation Research Board, 2010.
- [13] F. L. Mannering and S. S. Washburn, *Principles of Highway Engineering and Traffic Analysis*. Wiley, Third Edition, 2004.
- [14] "UPS," <http://www.ups.com/>.
- [15] "Seagate," <http://www.seagate.com/fr/fr/internal-hard-drives/enterprise-hard-drives/>.
- [16] M. Grünig, M. Witte, D. Marcellino, J. Selig, and H. van Essen, *Impact of Electric Vehicles: An overview of Electric Vehicles on the market and in development*. European Commission, 2011.
- [17] "Électricité de France," <http://particuliers.edf.com/abonnement-et-contrat/les-prix/les-prix-de-l-electricite/tarif-bleu-47798.html>.
- [18] "FranceIX," https://www.franceix.net/page.php?MP=editorial&ST=section&op=aff_section&secid=6.
- [19] "Better Place," <http://france.betterplace.com/>.
- [20] "Nanyang technological university," http://news.ntu.edu.sg/pages/newsdetail.aspx?URL=http%3A%2F%2Fnews.ntu.edu.sg%2Fnews%2FPages%2FNR2012_May24.aspx&Guid=b0330082-98bd-467b-97b9-f335d002aed7&Category=All/.
- [21] M. Al-Fares, A. Loukissas, and A. Vahdat, "A scalable, commodity data center network architecture," in *ACM Sigcomm*, Seattle, Washington, USA, Aug. 2008.
- [22] "Sandvine," <http://www.sandvine.com/>.
- [23] W. Zhao, Y. Chen, M. Ammar, M. Corner, B. Levine, and E. Zegura, "Capacity enhancement using throwboxes in dtms," in *IEEE International Conference on Mobile Ad-hoc and Sensor Systems*, Vancouver, Canada, Oct. 2006.
- [24] R. C. Shah, S. Roy, S. Jain, and W. Brunette, "Data mules: Modeling and analysis of a three-tier architecture for sparse sensor networks," *Ad Hoc Networks*, vol. 1, no. 2-3, pp. 215–233, Sep. 2003.
- [25] W. Zhao and M. H. Ammar, "Message ferrying: proactive routing in highly-partitioned wireless ad hoc networks," in *Workshop on Future Trends of Distributed Computing Systems*, San Juan, Puerto Rico, May 2003.
- [26] A. Keränen and J. Ott, "Dtn over aerial carriers," in *Workshop on Challenged Networks*, Beijing, China, Sep. 2009.