

# The Impact of Consumer Multi-homing Behavior on Ad Prices: Evidence from an Online Marketplace

Yu-Hsin Liu\*

January 2, 2018

## Abstract

This study examines whether consumer multi-homing behavior affects ad prices in the digital display advertising market via unwanted duplications. Theoretical analyses based on such channel draw important implications on media content and market structure. The technological advancement on frequency capping however may eliminate the over impression concern at the expenses of consumer privacy. The presenting paper provides the first empirical evidence on the mechanism. I scrape the publisher data from BuySellAds and match it with comScore 2016 for consumer multi-homing behavior. I employ a novel difference-in-difference strategy to identify the multi-homing effect on ad prices. The idea is that more viewable ad locations in a webpage are also more vulnerable to consumer multi-homing. By finding that the marginal effect of multi-homing (treatment) on ad prices is indeed more negative for the more viewable ads (treated), I conclude that consumer multi-homing behavior can increase the tendency of over impressions, and such tendency can decrease advertisers' valuation of ad slots in the digital display ad market.

---

\*Indiana University, Department of Business Economics and Public Policy, Kelley School of Business. I deeply thank Jeffrey Prince, Michael Baye, Matthijs Wildenbeest, Weihua An, and Stefan Weiergraeber for their valuable inputs and suggestions. I am also appreciative of the comments received from the faculty and the doctoral students in the BEPP and Economics departments. I am responsible for all errors.

# 1 Introduction

The Kia commercial “Save the Australian Open with Kia” that teams up Rafael Nadal and X-Men in 2015 creates surprising disdain in the YouTube community. Despite the great lineup that is supposed to match the taste of its target audiences and despite over 13 million views, there are more “dislikes” than “likes”<sup>1</sup>. As one of the many negative comments states: *“This advert makes me hate YouTube. It even makes me hate myself sometimes. Why would you start an advert with loud disturbing screeching sounds and pay for it to play before every video on YouTube?”* An industry commentator indicates that Kia may just skip the “frequency capping” function that limits the campaign impressions to each unique user<sup>2</sup>. The result is that the campaign budget burns out quickly with little and over-exposed reach. The careless frequency management seems to ruin the potential of the campaign.

Although frequency capping may benefit both advertisers and consumers, the challenge is on the behavioral tracking across websites. In the digital ad market<sup>3</sup>, frequency capping is implementable by storing parts of the consumer web browsing information in the ad server. For the frequency capping to apply to different media (e.g. Facebook, YouTube) or ad networks (e.g. AdWords), it requires that different ad servers share the browsing information. The sharing of the behavioral data can raise a substantial privacy concern. In the Federal Trade Commission (FTC) report of self-regulatory principles for online behavioral advertising (FTC, 2009), it states *“first-party (or intra-site) behavioral advertising practices are more likely to be consistent with consumer expectations, and less likely to lead to consumer harm,”* and yet *“when behavioral advertising involves the sharing of data with ad networks or other third parties, the consumer may not understand why he has received ads from unknown marketers based on his activities at an assortment of previously visited websites.”* Based on the

---

<sup>1</sup>[https://www.youtube.com/watch?time\\_continue=62&v=j2jm4r67Hjs](https://www.youtube.com/watch?time_continue=62&v=j2jm4r67Hjs), accessed in September 2017.

<sup>2</sup><https://dailycommercials.com/kia-kia-x-men/>, accessed in September 2017.

<sup>3</sup>The US digital ad market has a total revenue \$72 billion in 2016 and includes display ads (48%), search ads (46%), and others (6%, e.g. emails). The display ads include banner (40%), video (30%, e.g. YouTube commercials), rich media (20%, e.g. expendables), and sponsorship (10%, e.g. embedded logos) according to eMarketer (link).

argument, the “third-party” data sharing is excluded from the market practice, though the cross-website tracking is allowed within the same ad network. Given the market practice, the frequency management can still be challenging, especially when consumers tend to patronize different sets of online publishers (i.e. the consumer multi-homing behavior) in the digital space. That is, even if Kia uses the capping function in YouTube, their target audiences may still be over-exposed from seeing the commercial on other websites. In this paper, I empirically examine whether consumer multi-homing behavior affects ad prices via unwanted duplications in the online display ad market. The answer can reveal an important channel on how consumer privacy protection may affect the ad market (Budak et al., 2016; Goldfarb and Tucker, 2011; Johnson, 2013), and equally interestingly, it can have far-reaching implications on content provision in the consumer market (Athey et al., 2016; Anderson et al., 2015; Gentzkow and Shapiro, 2011).

Consumer multi-homing behavior is not specific to the digital space. Preceding papers recognize its role and importance in the broadcast industry (Anderson et al., 2015), the newspaper industry (Gentzkow et al., 2014), and the print magazine industry (Shi, 2015). They argue that ignoring this fact in such two-sided markets can lead to unrealistic market outcome predictions in the cases of media entry or merger and acquisition (Ambrus et al., 2014; Anderson et al., 2015; Anderson and Jullien, 2015). Some consequences of consumer multi-homing behavior can be far-reaching, as media publishers can alter contents to receive more valuable attentions from exclusive readers (Anderson et al., 2015) or to achieve more reach by expanding content breadth at the expense of content depth (Athey et al., 2016). In the political context, the pursuit of exclusive readers can arguably raise the issue of ideological segregation (Gentzkow and Shapiro, 2011; Gentzkow et al., 2014). The common implication is that a multi-homer’s attention is less valuable than a single-homer’s because of the potential ad duplication, and this incentivizes the economic agents and alters market outcomes. The empirical evidence in the offline market from Gentzkow et al. (2014) and

Shi (2015) supports the theoretical prediction of lower multi-homer value<sup>4</sup>. Given that the fundamental principle of the online media market is not very different from the offline, it is both natural and important to ask if the analyses in the offline context can apply to the digital display ad market – a \$35 billion business that generates 1.2 times the ad revenue of newspapers and print magazines combined<sup>5</sup>. To the best of my knowledge, this study is the first attempt to empirically examine the consumer multi-homing effect on ad prices in the online media market.

This paper addresses two important questions. The first question is whether consumer multi-homing behavior leads to over impressions. One specific point of examination is on the practice of frequency capping across websites. As the cross-website tracking is allowed within the same ad network, it is expected that the duplication problem is lessened. If the duplication problem is not a concern anymore, further relaxing the privacy requirement on “third-party” information sharing on consumer browsing behavior may not bring more benefit to frequency management. The flip side of the argument is that it indicates that ad networks may possess a significant amount of behavioral data that is enough to avoid the duplication inefficiency, which raises a further privacy concern on the FTC self-regulatory principles. Another point of examination is on the relevant measures of multi-homing, as it can be defined based on narrow genre market, broad genre market, and the entire web space. These points are unique online and can potentially affect the implications on the competition landscape of websites. The second question is whether the tendency of having impressions beyond the frequency cap can significantly lower advertisers’ willingness-to-pay for ad slots. If duplication is not a concern for online advertisers, then consumer multi-homing behavior would lose its channel to impact the ad prices. For the existing economic analyses to apply to the digital display ad market, both parts need to have significant effects.

To address the research questions, I obtain data from two sources, comScore 2016 and

---

<sup>4</sup>Chandra and Kaiser (2014) find the opposite, but their context is on cross-channel multi-homing (i.e. print magazine and companion website), rather than cross-publisher.

<sup>5</sup>The source is from eMarketer (link) estimated in March 2016, accessed in September 2017.

BuySellAds marketplace. For the consumer information, I utilize the web browsing tracking information in comScore to construct the multi-homing measures for each user<sup>6</sup>. For the publisher information, I scrape the ad prices and other relevant variables from BuySellAds. I then match the two data sets, calculate the multi-homing level for publishers by using a weighted average, and form a publisher-month panel. The procedure that the (only) preceding papers (Gentzkow et al., 2014; Shi, 2015) take is to collect the publisher-level data, infer the consumer multi-homing behavior from the supplementary survey data (e.g. preference rank on publishers) with structural models, and construct an aggregate demand system in the advertising market. One particular concern for the structural model in the online context is the potential scale of multi-homing, and the scale either requires the estimation of a large set of parameters or requires more restrictive assumptions. Since consumer multi-homing behavior is directly observed, I propose to use a reduced-form approach.

Recognizing endogeneity concerns, I employ a difference-in-difference strategy to identify how the multi-homing measure affects advertisers' valuation on ad slots. The central identification idea is that more viewable ad locations in a webpage are also more vulnerable to consumer multi-homing as the viewability can translate to the duplication potential. I regard the ads located on the top of the webpage as the treated group, the remaining ads as the control group (or "less treated" group), the multi-homing level as a continuous treatment, and the ad prices as the outcome variable. Although the marginal effect of multi-homing estimated within each group has an endogeneity problem, due to the unobserved publisher characteristics, strategic pricing behaviors, and the advertiser selection, I argue that the difference of the marginal effect only leaves the advertiser selection as a concern, and it is tractable. With plausible assumptions for interpretation, my result implies that higher consumer multi-homing level, defined based on the narrow genre market, increases the tendency of over impressions, and the tendency of over impressions decreases advertisers' valuation of ad slots. The result is robust to sample (publisher) selection, various panel regression

---

<sup>6</sup>The comScore 2016 data contains the web browsing information for each "machine" id. I use "user" to replace "machine" in the article.

models, and instrumental variable estimation.

The implications are two-fold. First, the empirical result confirms that under the current privacy principles the frequency management across websites remains challenging. Given that cross-website tracking is allowed within an ad network, further expanding of the network size or merger and acquisition of ad networks can improve the accuracy of frequency management. Such improvement may however evoke a re-examination on the current practice of consumer privacy protection. As previous papers (Goldfarb and Tucker, 2011; Johnson, 2013) indicate the economic trade-off between privacy protection and online ad effectiveness, a back-of-the-envelope calculation based on my estimation results shows that if cross-website tracking is perfectly allowed, the total revenue of the online display ad market can increase by at least 7%, and this only reflects the improvement on frequency management. Second, as the results indicate that the duplication problem is indeed a concern in a narrowly-defined genre market, there can be two potential content strategies for publishers to increase ad revenue. On the one hand, publishers may want to attract single-homing or “loyal” readers in a given genre, so they tend to explore niche markets. By having (more) single-homing readers, publishers can avoid the duplication discount. On the other hand, publishers may want to prevent advertisers from using other competing websites, so they may produce similar content to other competitors to widen the user base. By having (more) single-homing advertisers, publishers can also avoid the duplication discount. The first strategy echos (Anderson et al., 2015) and the second strategy echos (Athey et al., 2016). In addition to the changing ad revenue, privacy protection may also incentivize the use of these strategies.

I arrange the rest of the paper as follows. In Section 2, I provide an overall review of the economic analyses on how consumer multi-homing pattern may impact the media market. In Section 3, I provide a brief introduction on the current development of the digital advertising market. I then detail the data collection and the descriptive statistics in Section 4 and my empirical strategy in Section 5. In section 6, I report the results and discuss the implications. Section 7 provides the result from auxiliary models. Section 8 concludes.

## 2 Economic Theory Review

Consumer multi-homing behavior, where consumers tend to patronize multiple competing platforms for a similar service, is prominent in many modern platform businesses. To name a few, on-demand ride-sharing passengers may join more than one App platform (e.g. Uber, Lyft, etc); credit card holders may use more than one payment system (e.g. VISA, AmEx, etc); and newspaper readers may consume information from more than one source (e.g. New York Times, Washington Post, etc.). In many industries, this behavior may be reasonably regarded as a dynamic platform competition where the consumer switching cost is nearly zero, and the static single-homing model can provide a good basis for economic analyses (Armstrong, 2006; White and Weyl, 2016).

The media industry is however different. While each (repeated) match of the users from each side generate the same network effect in the ride-sharing or credit card platforms, in the media industry duplicated ad impressions may be wasteful. Theoretically, I see two major streams on analyzing how consumer multi-homing behavior influences the media market. Both streams post the assumptions of wasteful duplication and imperfect cross-publisher frequency management, which I empirically examine in this paper.

### 2.1 Incremental Pricing Principle

The first key channel is the Incremental Pricing Principle (Anderson et al., 2015; Anderson and Jullien, 2015) : given the diminishing return of the ad impressions, the optimal pricing rule for publishers indicates a lower ad price for a higher multi-homer composition. In their prototype model, the unique Nash equilibrium is that publishers implement the incremental pricing principle and advertisers patronize all publishers. The clean result comes from several strong, simplifying assumptions. Essentially, they rule out the advertiser heterogeneity and consumer strategic behaviors, purifying the multi-homing effect on ad prices to simply the average duplication discount.

Their motivation of relaxing the standard “single-homing” assumption is empirically driven. Earlier analyses based on two-sided markets (Anderson and Coate, 2015; Rysman, 2004) imply that the direct competition takes place on the consumer market and the ad slots of different publishers are not substitutes. The platform operators hence hold the monopoly position to deliver exclusive attentions to advertisers. When the media is purely ad-financed and given that consumers dislike ads (which is common in broadcast/TV as ads are more intrusive), the implication is that more (less) intensive competition on the consumer market leads to lower (higher) ad levels (regarded as the “prices” for consumer), and thus higher (lower) ad prices in equilibrium.

Several empirical findings, however, do not conform with the prediction. For example, in the 1996-2006 merger wave in the US radio industry, Jeziorski (2014) finds that ad prices in fact increase as a result of the merger, which even overcomes the potential price premium on audience homogeneity on niche markets Chandra (2009). In cases of deregulating public broadcasters to air ads, the single-homing model would predict that private, ad-financed broadcasters may be benefited by stealing some exclusive audiences those who would prefer private broadcasters if ads were allowed in the public broadcasters from the public stations. Nevertheless, Sthmeier and Wenzel (2012) documents that lobbying presented against the policy in several European markets, which contradicts the prediction based on the single-homing consumer assumption. By adapting the multi-homing set-up, (Anderson et al., 2015) show that in the case of broadcaster entry (merge), ad prices go down (up) through the potentially increasing (decreasing) composition of multi-homer, which connects to the empirical findings more tightly.

In addition to market structure and regulation, consumer multi-homing behavior can have far-reaching effects on the media content provision. The specific effect depends on the channel. When the incremental pricing principle applies, the multi-homing discount incentivizes publishers to pursue exclusive readers. In their Hotelling, two-publisher setting that allows readers to patronize both, in equilibrium the two publishers would fall apart in



the content spectrum, which is different from the standard case with consumer single-homing where both publishers locate themselves in the middle. By locating differently, each of them owns some exclusive readers and preserve higher ad prices. Although a general  $n$ -publisher setting should be extended with more cautions, the essential message is that when consumer multi-home, niche markets may receive too many content provisions from publishers while the general market receives too few. This contrasts with the early work (Steiner, 1952) that concludes that market equilibrium can lead to too many duplicated popular content provisions. In the political context, a potential concern driven by the incentive to capture exclusive attentions is ideological segregation (Gentzkow and Shapiro, 2011; Gentzkow et al., 2014). Interestingly, this is not rooted from consumer selection on the media channel (the so-called echo chamber as in Flexman et al. (2016)), but from publisher content strategies to attract more exclusive attentions.

## 2.2 Greater-reach Premium

The second key channel is the “greater reach premium.” Athey et al. (2016) derive the advertiser “mixed-homing” outcome and specify the second crucial stream with the consideration of the online media market. Their two-publisher model allows advertisers to have different valuations to each consumer attention. They show that in equilibrium, low-type advertisers choose to single-home on the publisher with greater reach to avoid any wasteful duplication, whereas high-type advertisers multi-home both publishers as they can tolerance some of the efficiency loss from having duplicated impressions<sup>7</sup>. Therefore, the publisher with superior reachability is prioritized by the advertisers, while the reduced advertising demand depresses the per-consumer ad price of the “less-reach” publisher. In the case of two publishers, this exactly conforms with the composition of multi-homers – the “less-reach” publisher must have higher multi-homing composition. Their model, therefore, links the marketing concept of “greater-reach” premium (per consumer) to the consumer multi-homing effect. The

---

<sup>7</sup>Athey et al. (2016) assumes that duplicated impression can be avoid within a publisher, but not across publishers, due to the imperfect cross-website tracking.

two concepts, however, may fall apart when there are more than two publishers. See more discussions on Section 7.1.

Interestingly, the pursuing of the greater-reach premium implies that publishers may attempt to widen their contents, which is opposite to the niche market strategy discussed in the previous section. Given the resource constraint, publishers often choose to produce between wild but shallow content or deep, full content. When the publisher with greater reach receives the priority on advertising demand, they show that the market equilibrium skews to the wilder but shallower content. The welfare implication is then the concern on the depth of online content, as it has become increasingly popular over time.

### 3 Industry Overview

With the advancement of sophisticated technologies, the dynamic of this industry can indeed be different from its offline counterpart. In this section, I provide an overview of the digital advertising industry, and I argue that the current practice in the digital display ad market does not alter the analytical foundation from the economic theory.

The US digital ad market has tripled its size from 2009 (\$23 billion) to 2016 (\$72 billion)<sup>8</sup>. In 2016, the major contributors are display advertising and search advertising, accounting for 48% and 46% of the total ad revenue respectively. Within the digital display ad market, eMarketer classifies the ad formats into four categories: banners (42%, including Facebook News Feed), video (30%, e.g. YouTube commercials), rich media (23%, e.g. expendables), and sponsorships (5%, e.g. embedded logos)<sup>9</sup>. Video and rich media are predicted to grow relatively faster in the next two years<sup>10</sup>. It is reasonable that the economic theory on the multi-homing effect should apply to all four categories, as in principle they are similar to offline media (e.g. newspapers, magazines, broadcast, etc.). Note however that the analysis may not be over-generalized to search advertising for two reasons. First, the multi-homing

---

<sup>8</sup>Estimated by the Interactive Advertising Bureau ([link](#)).

<sup>9</sup>Estimated by eMarketer ([link](#)).

<sup>10</sup>Estimated by eMarketer ([link](#)). See slide 13.

discount comes from the (potential) duplicated ad impressions that are less valuable, yet the duplication generated by user’s intent for multiple searches may actually mean the opposite. Second, unlike the display ad market that often uses the flat rate (e.g. per month) or cost per thousand impressions (cpm), the common pricing unit in the search ad market is per click. It is not clear that how multiple searches and impressions may translate to higher or lower price per click.

In the digital display ad market, publishers often sell their ad inventories through multiple prioritized sources. The basic idea is that publishers would sell their premium ad inventories in a relatively direct manner and clear the remnant inventories as much as they can with the assistance of ad technology. Whenever possible, publishers sell the ad inventories directly to advertisers in a negotiable and usually premium price. When direct sales are not possible or too costly, publishers follow similar logic to use the service from ad networks to monetize the ad slots.

Digital ad networks aggregate the ad supply from publishers and match it with advertiser demand using programmatic tools. The transaction types and the corresponding technology that the ad networks use fall into two categories: **programmatic direct** (56%) and **real-time bidding RTB** (44%)<sup>11</sup>. According to the Interactive Advertising Bureau (Interactive Advertising Bureau, 2013), publishers usually sell their premium and foreseeable ad inventories by adapting the programmatic direct service. They set or negotiate an up-front fixed rate for these ad slots, with the slots being guaranteed (**automated guaranteed transactions**) or prioritized (**preferred deals**) to sell to the buyers. This is similar to direct sales, except that a middleman implements the programmatic technology to facilitate ad inventory transactions and fill-ups. As the remaining ad supply becomes less predictable, publishers can adapt the real-time bidding service and set a **private auction** or an **open auction** to clear the remnant inventories. For the RTB services, advertisers usually do not know where their ads will be shown – they bid for the impression of some certain targeted audiences. A

---

<sup>11</sup>Estimated by eMarketer in April 2017 (link), accessed September 2017.

few ad networks that provide the services include: Google AdSense, the primary ad network that implements the open auction; Google Ad Exchange, which further includes the features of private auction and preferred deals, allowing more transparent information between ad buyers and sellers (e.g. publisher brand)<sup>12</sup>; and BuySellAds, where my publisher data comes from, operates automated guaranteed transactions and is a market leader of its kind.

## 4 Data

I collect the publisher data from BuySellAds ([www.buysellads.com](http://www.buysellads.com)), enrich it with the web categorizer from SimilarWeb ([www.similarweb.com](http://www.similarweb.com)), and match it with the consumer-level data from comScore 2016. To do so, I scrape the historical data in 2016 from the Internet Wayback Machine Archive ([www.web.archive.org](http://www.web.archive.org)).

### 4.1 BuySellAds

The publisher-level data comes from the BuySellAds marketplace. Although the financial reports for such market is often not publicly available, Business Wires claims that BuySellAds is “the first mover in Automated Guaranteed and market leader by revenue and transaction volume.”<sup>13</sup> It provides programmatic tools to clients, such as Bitcoin, Lonely Planet, Roku, Stack Exchange, etc<sup>14</sup>. It also operates an online marketplace for publishers to list their ad inventories and transact directly with advertisers. Some example publishers include The Atlantic, National Public Radio ([npr.org](http://npr.org)), Reddit, TechCrunch, etc.

The BuySellAds marketplace creates a transparent transaction environment that is similar to direct sales. Qualified publishers list the basic information of their ad inventories, including the website introduction, the display format, size, and location, with a non-negotiable

---

<sup>12</sup>See Google Support. Also, a rising concern of using the RTB service is its lack of transparency and controllability, and having ads on a fake news site, for example, is certainly not desired. As Google predicts that there is a rising trend of programmatic direct (link), they also adapt a more transparent design in the Ad Exchange and remove “bad publishers” from the AdSense network.

<sup>13</sup>See Yahoo Finance, reported November 2016, accessed September 2017.

<sup>14</sup>See <https://www.buysellads.com/about#timeline>.

posting price either in cpm or in a flat price over a month. This information essentially provides me the possibility of implementing the difference-in-difference strategy as different prices for different webpage locations and other ad specs are observed. BuySellAds categorizes publishers according to their content topics (e.g. Automotive, Entertainment, Technology, etc.). The categorization provides useful information on identifying the relevant genre market when developing the consumer multi-homing measures, as I will discuss more in Section 4.3. In addition to its middleman role, BuySellAds provides estimated impressions for any listing ad and allows both content targeting (e.g. ad paired with a certain content topic) and geographical targeting (e.g. US only). Advertisers can then shop those ad inventories in the same way that consumers shop online. See Figure 2 for the screen shots of the marketplace listings and Figure 3 for the screen shots of the buyer interface and the targeting options.

There can be, however, two concerns of using the publisher data from BuySellAds. First, the service fee that BuySellAds charges for listing ad inventories is not observed. It is possible that BuySellAds price discriminates different publishers and hence affects the observed listing prices of ad inventories. Second, publishers' posting prices can involve strategic decisions. As mentioned in Section 3, publishers often prioritize their ad selling channels. It is hence expected that their listing prices in the automated guaranteed marketplace can be strategically higher, backed up by other ad network services to clear up the remnant inventories, likely through Real-Time Bidding<sup>15</sup>. Although both concerns may raise the observed ad prices, there is no obvious reason on how the rising pattern can be affected by the publisher multi-homing level. If the price is just systematically higher, it does not affect the

---

<sup>15</sup>In fact, this back-up strategy may be the channel through which publishers learn the effective cpm (ecpm) its ad has. One may ask how publishers or advertisers may possibly know the consumer multi-homing pattern. One way is through the third-party service, such as Alexa or SimilarWeb who reports the overlapping measures between websites. The other mechanism is that publishers may not know the information directly, but the multi-homing pattern and the duplication effect have impacts on the Real-Time-Bidding service. When the multi-homing level is severe, frequency capping in the RTB service may prevent such publisher from receiving a higher bid because its reader may already be exposed with some ad impressions. The ecpm is hence lower. As this strategic reference point is lower, it may be optimal for publishers to set the listing price in BuySellAds lower too.

identification of the multi-homing effect.

To match the comScore 2016 data on the consumer side, I scrape the historical BuySellAds data from the Internet Wayback Machine Archive. The nonprofit organization archives over 305 billion webpages from 1996 automatically by its web crawler. One challenge of collecting data from the organization is that it only preserves partial information of a website on random dates. For example, when accessing the BuySellAds marketplace homepage on January 1st, the hyperlink to, say, the Business & Finance category from that homepage is redirected to March 19th, the nearest date that the Business & Finance webpage is preserved. As the web directory goes deeper, the number of preserved historical webpages becomes smaller, as shown in Figure 4. The data structure is thus a short, highly imbalanced panel<sup>16</sup>. The key variables collected include the ad prices, sizes, shape, and locations. The summary statistics are reported in Table 1.

## 4.2 ComScore 2016

The consumer-level data comes from comScore 2016. The data contains the Internet browsing history and demographic information for over 80,000 anonymous US households with their consent. The comScore data uniquely offers the crucial information I need for developing multi-homing measures and to retrieve some necessary but missing information in the advertising market. Specifically, the web-browsing history allows me to observe consumer multi-homing behavior precisely. Such data also allows me to calculate the market share of each website in its genre market and the average time duration consumers spend on each website<sup>17</sup>. I use the information for inferring the unobserved content quality and advertising level later in my empirical model to identify the multi-homing effect.

One major concern for matching the BuySellAds data with the comScore data is that

---

<sup>16</sup>The publisher data is concentrated around three particular dates: March 19th and April 10th in 2016, and January 7th in 2017.

<sup>17</sup>A potential alternative for (2) is publisher-level data on “average visiting time” from other web analytic services, such as Alexa. Alexa offers “visitor loyalty metrics” that provide estimates on page views per visit and per visitors. These can be a good proxy for duration, but coarse.

Table 1: Cpm Prices by Ad Specs

	Mean	S.D.	Min.	Max.	Obs.
<b>Location</b>					
top left	5.51	10.67	0.25	90	88
top center	6.68	6.23	0.10	35	321
top right	6.58	6.77	0.27	40	389
middle left	3.11	2.73	0.50	12	38
middle center	4.75	4.62	0.25	20	96
middle right	4.33	5.41	0.25	35	156
bottom left	2.97	4.00	0.10	15	24
bottom center	3.16	3.32	0.25	16	85
bottom right	3.21	3.87	0.30	18	54
<b>Size*</b>					
large (area $> 75000$ dips)	9.16	7.88	0.25	40	207
medium	4.98	6.53	0.10	90	585
small (area $\leq 65520$ dips)	4.63	4.74	0.10	30	464
<b>Shape</b>					
banner	5.49	5.68	0.10	30	494
square	4.90	6.41	0.10	90	612
skyscraper	8.32	7.76	0.25	40	150
<b>Unknown</b>	4.26	4.33	0.10	35	843

\*The size of digital ads are measured in density-independent pixels (dips). There are some common sizes for ads standardized by the Interactive Advertising Bureau (IAB) but in principle they can be any size. Some common sizes are  $728 \times 90$  (small banner),  $125 \times 125$  (small square),  $300 \times 250$  (medium square), and  $300 \times 600$  (large skyscraper).

despite the fact that there are over 80,000 households in comScore, publishers in BuySellAds may receive only a few visits from the observed samples. This is not surprised considering that there are more than 120 million households in the US. To develop the measure of multi-homing and other the non-price characteristics of the publishers, such as the average reader profile on income, education, duration, I however rely on those observed samples. A compromise is thus made to take the data panel to a monthly level to increase the representativeness of the measures of the publisher features. In my main model, I include the publishers with at least 20 unique observed users in a given month to avoid noisy measures. I also do a robustness check on setting different thresholds.

### 4.3 Measures for “Multi-homing” and SimilarWeb

I now describe the measures of “multi-homing” for online publishers. I use a weighted average (over visitors) of the multi-homing numbers as the continuous measure of publisher multi-homing level. In my data, no online publisher owns exclusive attentions. Hence a dichotomy measure similar to Gentzkow et al. (2014) and Shi (2015) is not proper in this context. To construct the measure, I use the web-browsing tracking data from the comScore 2016. I first count the number of websites a household visited on a daily basis, as a daily ad cap is the most common market practice<sup>18</sup>. For each publisher, I then take a weighted average over its visitors in each month with the weight being the number of its webpages consumed by that visitor. For example, suppose there is a website browsed by two visitors, one is a 5-homer in that genre and consumes one webpage of the website; another is a 3-homer in that genre and consumes four webpages of that website. I give 20% to the 5-homer and 80% to the 3-homer. The multi-homing level for the website is then 3.4. The weighting makes sense as it reflects the actual likelihood that an ad in that website impresses the corresponding multi-homing household.

---

<sup>18</sup>With the help of advanced tracking technology, online advertisers can put a daily cap and a weekly (monthly) cap for ad impressions by the same viewer, yet it is not reasonable to put only a weekly cap without a daily cap as it would be wasteful if the weekly cap is used up in one day. For the duplication discount to be reasonable, the daily measure is the most relevant measure.



The next question is then whether the multi-homing measure should count all websites, even if they have very different contents. Online advertisers may target a specific consumer in all kinds of websites and show a specific ad. They may also only want to put the ad on some websites in a specific genre for content relevance. If all websites are possible channels, it means that the duplication effect would count across all websites. If only the websites within a genre matter to advertisers, then the duplication effect should only apply to the specific genre. As there is no clear guide on what advertisers usually do and no previous research examines such a question, I develop three levels of multi-homing measures and test to which levels the duplication effect remains relevant.

To construct the multi-homing measures in three levels, I first need to categorize as many websites as I can. To do so, I scrape the web categorizer from SimilarWeb. The SimilarWeb free version provides the top 50 websites (in terms of worldwide web traffic) in each of the 169 sub-categories. I scrape the information and match it with the twenty categories that BuySellAds uses. In the first level, I use the twenty categories and regard each as a genre market to calculate the multi-homing measure for a publisher ( $H^{narrow}$ ). I measure the corresponding multi-homing behavior only within that genre. That is, if a website visitor visits another website that is not in the same genre, it does not count as multi-homing. In the second level, I group the twenty categories to eight broader genres and calculate the multi-homing level within the broad genre ( $H^{broad}$ ). Finally, I calculate the multi-homing level that counts every website ( $H^{universal}$ ).

The summary statistics for each of the twenty categories are reported in Table 2. For example, for the Business & Finance category, there are 46 publishers observed in the matched data. On average, the daily multi-homing level for a sample publisher is 1.88 out of the total 435 websites categorized to Business & Finance and 1.95 out of the broader Business & Politics genre. The measures can increase dramatically when using week or month as the multi-homing time window, indicating that the observed users do not typically consume the information from few fixed sources. Also, from the statistics of  $H^{universal}$ , there is no signifi-

Table 2: (Daily) Multi-homing Level of Sample Publishers by Genre

Category	Publishers		$H^{narrow}$		$H^{broad}$		$H^{universal}$	
	Categorized	Sample	Mean	S.D.	Mean	S.D.	Mean	S.D.
<b>Business &amp; Politics</b>								
Business & Finance	435	46	1.88	1.02	1.95	1.11	29.14	23.78
Government & Politics	219	22	1.14	0.18	1.61	0.39	24.40	8.92
<b>Entertainment</b>								
Entertainment	338	42	1.33	0.62	1.61	1.17	30.27	35.30
Gaming	823	15	1.92	0.75	2.11	0.72	15.61	8.25
<b>Fashion</b>								
Beauty & Fashion	260	10	1.20	0.39	1.26	0.57	32.46	48.30
Weddings	131	4	1.34	0.87	1.67	0.93	34.11	11.97
<b>Health</b>								
Food & Drink	798	10	1.84	0.81	1.88	0.80	25.61	13.33
Health & Fitness	260	10	1.15	0.20	1.34	0.50	29.89	22.22
<b>Life</b>								
Parenting & Education	357	37	1.32	0.37	1.36	0.40	21.65	9.97
Home & Architecture	340	12	1.48	0.64	1.66	0.66	36.42	29.51
Pets	874	5	1.55	0.67	1.78	0.63	20.93	5.22
<b>Recreation</b>								
Automotive	1057	15	2.01	1.11	2.15	1.29	16.14	7.14
Sports	2249	32	3.30	3.68	3.64	3.40	36.87	57.91
Travel	880	8	1.71	0.95	1.75	0.98	20.52	8.12
<b>Tech</b>								
All Things Apple	86	24	1.07	0.16	1.81	0.60	18.82	8.05
Cryptocurrency	13	1	1	0	2	0	23.5	0
Technology	807	71	2.29	0.88	2.35	0.91	28.13	24.17
Virtualization	26	4	1.03	0.13	2.46	1.66	69.05	174.76
<b>Web Design</b>								
Visual Arts & Design	146	31	1.06	0.17	1.87	1.04	27.78	13.10
Web Design & Development	512	130	2.34	1.52	2.41	1.55	29.67	24.63

cant different patterns between categories on how much a typical consumer may multi-home in general.

## 5 Empirical Strategy

### 5.1 Main Hypothesis

The research question is to examine if the consumer multi-homing behavior affects the advertising demand via the duplication effect. The multi-homing measures can regard the narrow, broad, or entire web space, but I do not distinguish them in this section for avoiding redundancy.

I now illustrate the empirical question and the identification strategy formally. Denote  $H_k$  as the total number of online publishers in the relevant content market that a visiting consumer  $k$  patronizes. Let the unobserved duplication level  $D(H_k)$  be the number of “over” impression(s) of an ad seen by the consumer  $k$  before visiting the current publisher. By “over” I mean the number of impressions that exceed the frequency cap that the advertiser has in mind but is unable to perfectly implement. Assume that  $D(\cdot) \geq 0$  and  $D'(\cdot) \geq 0$ .  $D'(\cdot) > 0$  only if both the following cases are true. First, advertisers utilize multiple publishers for reaching targeted audiences, so that the duplication across publishers has positive probability. Second, there is no universal frequency capping applied, otherwise  $D(\cdot)$  will always be zero. For simplicity, I further assume that  $D(\cdot)$  is a linear function. Hence, the expected duplication level is  $E[D(H_k)] = D[E(H_k)] = D(H)$ .

For advertisers, I assume that the value of an impression  $V$  in each publisher depends on its expected duplication level  $D$ . The definition of  $D(\cdot)$  distinguishes my study from other studies that concern the optimal number of impressions, where the marginal value of an impression is commonly found to first increase and then decrease. In the current setup, since it is not optimal for advertisers to set the frequency cap on where the marginal value is still rising, one more duplication over the frequency cap is expected to have a lower

incremental value. The main hypothesis is

$$V_H = \frac{\partial V}{\partial D} \frac{\partial D}{\partial H} < 0$$

If  $V_H < 0$ , it implies  $\frac{\partial D}{\partial H} > 0$  and  $\frac{\partial V}{\partial D} < 0$ . It further implies that the universal frequency capping is not the market practice and the consumer multi-homing behavior generates a per-impression price discount via the duplication effect.

## 5.2 Difference-in-difference Strategy

Ideally, to identify  $V_H$  it would be sufficient to vary  $H$  of a publisher and examine how its ad demand would shift. Practically, this is not possible as one cannot experimentally change  $H$  and the demand curve is not observed. What can be observed are the consumer multi-homing behaviors for different publishers and the market prices of the ads. When using such data to interpret  $V_H$ , the endogeneity problem can be severe. On the one hand, there are other reader characteristics that may correlate with  $H$  and affect advertiser valuation for an impression. On the other hand, the variation of market prices does not necessarily capture the equivalent information of the variation of advertiser valuation, due to the strategic and simultaneous natures of the market equilibrium. Hence, even with sufficient controls on reader characteristics,  $V_H$  is unlikely to be identified using only the cross-publisher variation.

One way to address the endogeneity problem is to construct a quasi-experiment. The identification idea is to find a group of ads for which the duplication level is more responsive to the multi-homing level and this is the only difference between it and the other group of ads. Next, one can observe the variation of the multi-homing level within each group and estimate its marginal effect on ad prices. The marginal effect within each group cannot tell a causal story because of the endogeneity problem. However, by comparing the marginal effect between the two groups, one can reasonably claim that the difference of the marginal effect is due to the heterogeneous duplication effect. In practice, such almost identical treated

and control groups are difficult to find, but for identification the assumption can be further relaxed. The essential assumption is that the treated and the control can be different, but the unobserved heterogeneity does not correlate with both the ad prices and the publisher multi-homing level. Hence, for searching the quasi-experimental potential, the criteria is to (1) find the treated and control groups (2) that are otherwise as similar as possible, and (3) for the remaining unobserved heterogeneity between the groups, rely on the plausible “uncorrelation” assumptions in the given context.

The grouping variable I use for the quasi-experiment is the ad location on a webpage, and I argue that it satisfies the aforementioned criteria. First, the ads located on the “top” of the webpage are more viewable than others, and thus they should also be more responsive to consumer multi-homing behaviors. That is, denote  $T$  as the top and treated ad group, and  $B$  as the non-top and control group,  $\frac{\partial D^T}{\partial H} \geq \frac{\partial D^B}{\partial H} \geq 0$ . For advertisers who seek for more viewable locations to post their ads, they are likely to purchase the “top” ads across publishers. In this case, a publisher with high  $H$  can be particularly concerned as the viewability can translate to more wasteful duplication. The opposite argument can be given for the less viewable locations. Second, In my data most publishers have both “T” ads and “B” ads. So for the ads in both groups, they are balanced in terms of both publisher and reader characteristics. Third, the remaining endogeneity problem comes from publisher strategic behavior on pricing “T” and “B” ads and also from the advertiser selection. For the publisher strategic behavior, I allow such behavior but assume that there is a general pricing strategy towards “T” and “B” ads that does not correlate with  $H$ . For the advertiser selection, I assume that the ads posted on the “T” location do not have systematic differences that correlate with  $H$ . The assumption on the advertiser selection can potentially be problematic, and I will discuss it shortly.

Given the identification assumptions, although  $V_H$  may still not be identified with the observed market data,  $\Delta V_H$  can be identified more confidently. I state the altered hypothesis

under the difference-in-difference strategy as follows:

$$\Delta V_H = V_H^T - V_H^B = \frac{\partial V}{\partial D} \left( \frac{\partial D^T}{\partial H} - \frac{\partial D^B}{\partial H} \right) < 0$$

If  $\Delta V_H < 0$ , it implies  $\frac{\partial D^T}{\partial H} - \frac{\partial D^B}{\partial H} > 0$  and  $\frac{\partial V}{\partial D} < 0$ . Essentially I utilize the heterogeneous duplication responsiveness to identify the desired effect. The implication is exactly the same as the main hypothesis.

Finally, I note that the advertiser selection can cause endogeneity concerns and affect the implications. The impacts are two-fold. For the heterogeneous  $D(\cdot)$ , knowing the superior viewability, advertisers who use “T” slots may use a more restrictive set of publishers, and thus  $\frac{\partial D^T}{\partial H}$  may become less responsive. Also, when more advertisers use both “T” and “B” ad slots,  $\frac{\partial D^T}{\partial H}$  and  $\frac{\partial D^B}{\partial H}$  can be closer in magnitude. Hence, advertiser selection may reduce the heterogeneous duplication effect. In the extreme case where  $\frac{\partial D^T}{\partial H} - \frac{\partial D^B}{\partial H}$  is close to zero, the identification strategy cannot work. For the ad content selection, advertisers may pair the ad contents that are more “lasting” with the “T” ad locations. That is, the value of the ads may preserve relatively higher value as the number of duplications (over the frequency cap) increases,  $\frac{\partial V^T}{\partial D} \geq \frac{\partial V^B}{\partial D}$ . This concern creates an ambiguity on the implication of the sign of  $\Delta V_H$  as it can be different from  $\frac{\partial V}{\partial D}$ .

To be clear, I explicitly state the three relaxed key assumptions and how to interpret the result of  $\Delta V_H$  based on them. The assumptions are:

1.  $\Delta V_H = \frac{\partial V^T}{\partial D} \frac{\partial D^T}{\partial H} - \frac{\partial V^B}{\partial D} \frac{\partial D^B}{\partial H}$  is identified (instead of  $\frac{\partial V}{\partial D} \left( \frac{\partial D^T}{\partial H} - \frac{\partial D^B}{\partial H} \right)$  is identified)
2.  $\frac{\partial D^T}{\partial H} \geq \frac{\partial D^B}{\partial H} > 0$  or  $\frac{\partial D^T}{\partial H} = \frac{\partial D^B}{\partial H} = 0$
3.  $\frac{\partial V^T}{\partial D} \geq \frac{\partial V^B}{\partial D}$  (instead of  $\frac{\partial V^T}{\partial D} = \frac{\partial V^B}{\partial D} = \frac{\partial V}{\partial D}$ )

Assumption 2 rules out the unreasonable case where the consumer multi-homing behavior only leads to duplication for the “T” ads but not the “B” ads. Assumption 3 reflects advertisers’ adjustment to the viewability difference. Now suppose that the empirical result

indicates  $\Delta V_H < 0$ , it implies that  $\frac{\partial V^T}{\partial D} / \frac{\partial V^B}{\partial D} < \frac{\partial D^B}{\partial H} / \frac{\partial D^T}{\partial H} < 1$ . Given assumption 2 and 3, the only possibility is  $\frac{\partial D^T}{\partial H} > \frac{\partial D^B}{\partial H} > 0$  and  $\frac{\partial V^B}{\partial D} < \frac{\partial V^T}{\partial D} < 0$ . In this case, the duplication effect can still be claimed. The intuition is that the advertiser selection works against the duplication effect, and yet the effect is strong enough to still show a significant pattern. The case where  $\Delta V_H \geq 0$  is however trickier, as there can be multiple explanations and the sign of both  $\frac{\partial V^T}{\partial D}$  and  $\frac{\partial V^B}{\partial D}$  cannot be confirmed.

### 5.3 Econometrics

To estimate  $\Delta V_H$ , I use the following specification for publisher  $i$ , ad slot  $j$ , and month  $t$ . The main goal is to identify  $\lambda$ , which is a vector of difference-in-difference parameters.

$$P_{ijt} = \beta_0 + \gamma T_{ijt} + f^n(H_{it}|\delta) + f^n(T_{ijt} \times H_{it}|\lambda) + S'_{ijt}\beta_1 + X'_{it}\beta_2 + \nu_i + \tau_t + \epsilon_{ijt}$$

where  $P_{ijt}$  is the cpm price,  $T_{ijt} = 1$  if being treated (top ad), and  $H_{it}$  is the multi-homing level. I denote  $f^n(H_{it}|\delta) = \delta_1 H_{it} + \delta_2 H_{it}^2 + \dots + \delta_n H_{it}^n$  and  $f^n(T_{ijt} \times H_{it}|\lambda) = \lambda_1(T_{ijt} \times H_{it}) + \lambda_2(T_{ijt} \times H_{it}^2) + \dots + \lambda_n(T_{ijt} \times H_{it}^n)$ . In my main model, I set  $n = 2$ . For the control variables,  $S_{ijt}$  is a vector of ad specs (e.g. size and shape), and  $X_{it}$  is a vector of publisher characteristics and representative reader demographics. The error terms include publisher-invariant term  $\nu_i$ , time-invariant term  $\tau_t$ , and other general shocks that potentially include publisher strategic behaviors advertiser strategic behaviors  $\epsilon_{ijt}$ . These error term may prevent the identification of  $\delta$ , but not  $\lambda$ . To match the treated ads and the controls ads more precisely, in my main estimation I also include the publisher fixed effect and month fixed effect. Finally, I take a log form for both prices and the multi-homing levels to normalize their highly skewed (to the left) and light-tailed distribution patterns.

## 6 Results

### 6.1 The Relevant Multi-homing Measure

This section focuses on the results and the interpretations on  $\frac{\partial D}{\partial H}$ . The results in Table 3 indicate that only the narrow multi-homing measure ( $H^{narrow}$ ) matters to the duplication effect. In this difference-in-difference specification with both publisher and month fixed effects, the key variables of interest are the interaction terms. Although all the multi-homing measures do not appear to have significant interaction effects in the linear specification ( $n = 1$ ), the narrow measure (and only it) shows a significant pattern in the quadratic specification ( $n = 2$ ). The comparison of the results from using the three measures is meaningful. Recall the three assumptions: (1)  $\Delta V_H = \frac{\partial V^T}{\partial D} \frac{\partial D^T}{\partial H} - \frac{\partial V^B}{\partial D} \frac{\partial D^B}{\partial H}$  is identified, (2)  $\frac{\partial D^T}{\partial H} \geq \frac{\partial D^B}{\partial H} > 0$  or  $\frac{\partial D^T}{\partial H} = \frac{\partial D^B}{\partial H} = 0$ , and (3)  $\frac{\partial V^T}{\partial D} \geq \frac{\partial V^B}{\partial D}$ . Since  $V(\cdot)$  is exactly the same for all three measures, the result difference comes from and only from  $D(\cdot)$ . Since the result for the narrow measure indicates that  $\frac{\partial V^T}{\partial D}$  and  $\frac{\partial V^B}{\partial D}$  cannot both be zero, unless the rare coincidence  $\frac{\partial V^T}{\partial D} / \frac{\partial V^B}{\partial D} = \frac{\partial D^B}{\partial H} / \frac{\partial D^T}{\partial H}$  happens to be true, the only plausible explanation for the insignificant results for the broad and universal measures is that they have insignificant impacts on the duplication level. That is, the results in Table 3 imply  $\frac{\partial D^T}{\partial H} = \frac{\partial D^B}{\partial H} = 0$  for the broad and universal multi-homing measures. A likely interpretation is that advertisers tend to implement the cross-publisher marketing campaigns within a narrow genre. Because the target audiences have zero probability to receive duplicated impressions from visiting websites outside the genre, the broader multi-homing pattern does not translate to the duplication problem.

The finding reveals the competition features of the online media market. In the online media market, one may argue that every website competes against each other as potentially they may receive the attention from any consumer and resell it to any advertiser. However, my finding indicates that specific advertisers may only use publishers in a specific genre to reach target audiences. Hence, whether publishers across genres compete against each



Table 3: Results for Narrow, Broad, and Universal Multi-homing Measures

	Narrow		Broad		Universal	
	$\ln(cpm)$	$\ln(cpm)$	$\ln(cpm)$	$\ln(cpm)$	$\ln(cpm)$	$\ln(cpm)$
$T$	0.151 (1.38)	0.331*** (3.09)	0.0727 (0.57)	0.227 (1.29)	-0.0686 (-0.12)	3.293 (1.63)
$\ln(H)$	0.195 (0.97)	0.855** (2.38)	0.0552 (0.30)	0.581 (1.39)	0.00970 (0.06)	3.188** (2.22)
$\ln(H)^2$		-0.389*** (-2.67)		-0.294* (-1.74)		-0.499** (-2.17)
$T \times \ln(H)$	0.00848 (0.05)	-0.782*** (-3.23)	0.130 (0.81)	-0.373 (-0.99)	0.0711 (0.40)	-2.099 (-1.62)
$T \times \ln(H)^2$		0.527*** (4.52)		0.304* (1.86)		0.347 (1.66)
ad_large	0.292** (2.44)	0.293** (2.46)	0.293** (2.45)	0.294** (2.48)	0.293** (2.45)	0.296** (2.53)
ad_small	-0.180 (-1.09)	-0.173 (-1.06)	-0.185 (-1.11)	-0.180 (-1.10)	-0.181 (-1.10)	-0.179 (-1.10)
ad_skycraper	-0.347*** (-3.04)	-0.340*** (-3.04)	-0.352*** (-3.05)	-0.346*** (-3.02)	-0.348*** (-3.05)	-0.342*** (-3.02)
ad_square	-0.219 (-1.27)	-0.209 (-1.24)	-0.224 (-1.30)	-0.219 (-1.28)	-0.220 (-1.28)	-0.216 (-1.27)
_cons	1.530*** (7.32)	1.363*** (6.21)	1.605*** (7.40)	1.434*** (5.52)	1.606*** (2.77)	-3.413 (-1.51)
Publisher FE	Yes	Yes	Yes	Yes	Yes	Yes
Month FE	Yes	Yes	Yes	Yes	Yes	Yes
User base	$\geq 20$	$\geq 20$	$\geq 20$	$\geq 20$	$\geq 20$	$\geq 20$
Publishers	75	75	75	75	75	75
Observations	360	360	360	360	360	360

$t$  statistics in parentheses

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

standard error clustered in the publisher level

other seems unclear, but publishers within a genre do. The reasons could be that the ads are more effective with compatible web content and that a specific genre provides specific targeting information. In a genre market, publishers may focus on different content to attract valuable readers with less multi-homing pattern (Anderson et al., 2015), or they may expand the content breadth in the genre to increase the reach (Athey et al., 2016). Beyond the genre market, such strategies may not hold. For example, a publisher may not have the incentive to attract a reader that is universally less multi-homing. Also, a tennis publisher may want to extend the content to soccer, rather than to other genre topics, such as movie for the “relevant” greater reach.

To confirm that the result of using  $H^{narrow}$  with quadratic specification ( $n = 2$ ) is robust, I check the results with different thresholds of the observed user base, as shown in Table 4. For example, when the threshold is three (the 1st column), it means that the analysis only includes the publishers with at least three unique observed users in the comScore data in a given month. The choice of the threshold has its trade-off. A higher threshold improves the representativeness of the multi-homing measures and the reader profiles of that publisher, but it also means only the thicker publishers are included in the analysis. The estimates on the interaction terms however do not show different patterns, meaning that the results are robust to the threshold selection.

## 6.2 The Duplication Effect

This section focuses on the results and the interpretations on  $\frac{\partial V}{\partial D}$ . I first interpret the estimated coefficients substantially. Specifically, I take the model with threshold 20 (the 3rd column in Table 4) as the main model. On average, the prices for the “top” location and “large” size, and “banner” shape are about 33%, 29% (compare to the middle size), and 34% (compare to the skyscraper shape) higher respectively. The difference however may not only be due to the viewability and size premiums, but also possibly reflect the advertiser selection and the publisher pricing strategy.

Table 4: Results for Publishers with Thinner and Thicker Observed Users

	(1) $\ln(cpm)$	(2) $\ln(cpm)$	(3) $\ln(cpm)$	(4) $\ln(cpm)$
$T$	0.442*** (4.35)	0.279** (2.56)	0.331*** (3.09)	0.322*** (2.83)
$\ln(H)$	0.396 (1.06)	0.592* (1.79)	0.855** (2.38)	0.924** (2.37)
$\ln(H)^2$	0.0722 (0.34)	-0.274** (-2.08)	-0.389*** (-2.67)	-0.426*** (-2.70)
$T \times \ln(H)$	-0.846*** (-3.20)	-0.456* (-1.80)	-0.782*** (-3.23)	-0.663*** (-2.80)
$T \times \ln(H)^2$	0.575*** (3.63)	0.343*** (2.67)	0.527*** (4.52)	0.451*** (4.23)
ad_large	0.133 (1.28)	0.303*** (2.75)	0.293** (2.46)	0.293** (2.42)
ad_small	-0.168 (-1.12)	-0.131 (-0.85)	-0.173 (-1.06)	-0.202 (-1.20)
ad_skycraper	-0.0959 (-0.68)	-0.284** (-2.55)	-0.340*** (-3.04)	-0.369*** (-3.16)
ad_square	-0.0814 (-0.52)	-0.0997 (-0.62)	-0.209 (-1.24)	-0.187 (-1.07)
_cons	0.998*** (4.75)	1.255*** (6.09)	1.363*** (6.21)	1.362*** (5.74)
Publisher FE	Yes	Yes	Yes	Yes
Month FE	Yes	Yes	Yes	Yes
User base	$\geq 3$	$\geq 10$	$\geq 20$	$\geq 30$
Observations	823	496	360	291
Publishers	209	113	75	62

$t$  statistics in parentheses

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

standard error clustered in the publisher level

The key estimates are the coefficients of the interaction terms. As argued in Section 5, the cross-publisher variation on multi-homing levels in each group may not be enough for identifying  $V_H$  due to various endogeneity issues. By applying the difference-in-difference strategy, however,  $\Delta V_H$  can be identified by the observed data given plausible assumptions. As advertiser selection can be a concern, the specific difference-in-difference effect the model identifies is  $\Delta V_H = \frac{\partial V^T}{\partial D} \frac{\partial D^T}{\partial H} - \frac{\partial V^B}{\partial D} \frac{\partial D^B}{\partial H}$ . The estimated  $\Delta V_H$  is  $-0.782 + 1.054 \ln(H)$ . This means that when the multi-homing level increases by 1% from the single-homing baseline (i.e.  $\ln(H) = 0$ ), on average the cpm price of the “top” ad receives a more negative impact, which is about 0.78%. The exact meaning of the scale can be vague, as the scale captures both the heterogeneous duplication levels ( $\frac{\partial D^T}{\partial H}$  and  $\frac{\partial D^B}{\partial H}$ ) and heterogeneous duplication effects ( $\frac{\partial V^T}{\partial D}$  and  $\frac{\partial V^B}{\partial D}$ ). However, such information can help assess the differential economic impact of privacy regulation on “T” and “B” ads through frequency management, as I will discuss in the next sub-section.

As shown in Figure 1,  $\Delta V_H$  remains negative when  $\ln(H) < 0.5$  and becomes positive when  $\ln(H) > 1$ . Note also that the 50 percentile of  $\ln(H)$  is 0.48 and the 90 percentile is 1.06, indicating that most data points are concentrated in  $0 < \ln(H) < 0.5$ . Within this range,  $\Delta V_H < 0$ . Given the three key assumptions (1)  $\Delta V_H = \frac{\partial V^T}{\partial D} \frac{\partial D^T}{\partial H} - \frac{\partial V^B}{\partial D} \frac{\partial D^B}{\partial H}$  is identified, (2)  $\frac{\partial D^T}{\partial H} \geq \frac{\partial D^B}{\partial H} > 0$  or  $\frac{\partial D^T}{\partial H} = \frac{\partial D^B}{\partial H} = 0$ , and (3)  $\frac{\partial V^T}{\partial D} \geq \frac{\partial V^B}{\partial D}$ ,  $\Delta V_H < 0$  implies (1)  $\frac{\partial D^T}{\partial H} > \frac{\partial D^B}{\partial H} > 0$  and (2)  $\frac{\partial V^B}{\partial D} < \frac{\partial V^T}{\partial D} < 0$ .

The first implication says that there are heterogeneous duplication levels between the “T” and “B” ad groups caused by consumer multi-homing behavior. The duplication level of both “T” and “B” ads must be affected by the consumer multi-homing behavior. Since the duplication level is specifically defined as the impressions over the desired frequency cap, the result implies that the frequency cap does not perfectly apply across publishers. Hence, although the current self-regulation on privacy gives ad networks some advantage, it doesn’t completely solve the challenge on cross-web tracking.

The second implication says that the duplication is indeed a concern for online advertisers

and advertiser selection may work against the identification of the duplication effect. The duplication concern can be alleviated by carefully designing the ad content, and advertisers are reasonable to pair more lasting ads with more viewable locations. Hence, although the “T” ads may receive more duplications than the “B” ads in the same webpage, it does not necessary mean that the value of “T” ads is affected more by the multi-homing level. Advertiser selection however does not prevent the finding of the existence of the duplication discount. The duplication effect is strong enough to overcome the counter effect.

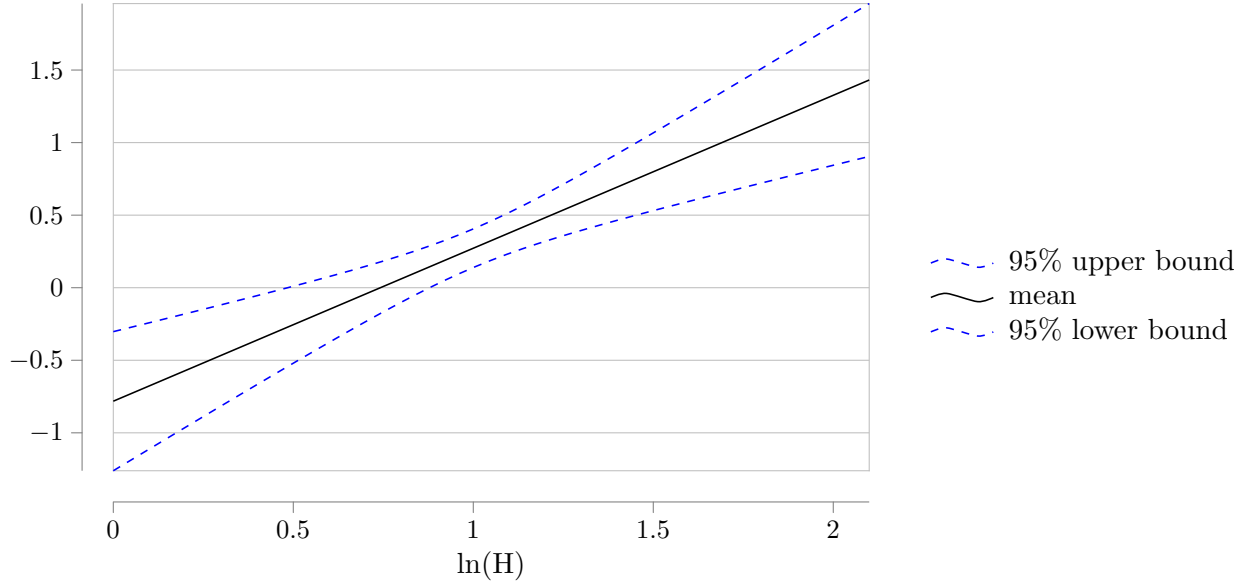


Figure 1: Estimated  $\Delta V_H$

Next, I explain the positive trend that leads  $\Delta V_H > 0$  when  $\ln(H) > 1$ . The broad idea is that I attribute the positive trend to advertiser selection, which works against the multi-homing effect. Specifically, I impose another assumption that  $V(D)$  is unimodal for both groups of ads. This means that once  $V_H$  becomes negative at a value ( $\bar{H}$ ),  $V_H \leq 0$  for all  $H \geq \bar{H}$ . That is to say, if the third duplication has less contribution than the second, the forth duplication cannot have more contribution than the third. I also recall the assumption that  $D(H)$  is linear for both groups. With the two augmented assumptions, the only explanation of  $\Delta V_H = \frac{\partial V^T}{\partial D} \frac{\partial D^T}{\partial H} - \frac{\partial V^B}{\partial D} \frac{\partial D^B}{\partial H}$  appearing a positive trend is the case where  $\frac{\partial V^B}{\partial D}$  becomes more negative relative to  $\frac{\partial V^T}{\partial D}$  as  $H$  and hence  $D$  increase.  $\Delta V_H$  flips to positive

when  $\frac{\partial V^T}{\partial D} / \frac{\partial V^B}{\partial D} < \frac{\partial D^B}{\partial H} / \frac{\partial D^T}{\partial H}$ . The intuition is that the ad content in the top locations can endure more duplication particularly when the potential duplication level is high. This makes sense as more lasting ad content and more viewable ad locations are strategic complements. Note that the third key assumption ( $\frac{\partial V^T}{\partial D} \geq \frac{\partial V^B}{\partial D}$ ) is based on the same “advertiser selection” argument, but the assumption is on the “first order” effect of the heterogeneity on the ad lasting feature. With the other assumptions, the estimating result suggests that the “second order” effect: not only  $\frac{\partial V^T}{\partial D} \geq \frac{\partial V^B}{\partial D}$  but the gap further becomes larger.

To summarize, given the three key assumptions,  $\Delta V_H < 0$  implies  $\frac{\partial V^B}{\partial D} < \frac{\partial V^T}{\partial D} < 0$ . Combined with the two augmented assumptions, the positive trend of  $\Delta V_H$  implies that  $\frac{\partial V^B}{\partial D} < \frac{\partial V^T}{\partial D} < 0$  and the gap between  $\frac{\partial V^B}{\partial D}$  and  $\frac{\partial V^T}{\partial D}$  increases as  $D$  increases. Since the empirical evidence supports that consumer multi-homing behavior indeed affects the advertising demand via the duplication discount, the existing economic analyses (Athey et al., 2016; Anderson et al., 2015) based on such channel can still matter for online media that are monetized through display ads under the current privacy law.

### 6.3 Privacy Protection and Economic Implications

The effectiveness of online behavioral advertising comes from detailed consumer behavioral data. The collection of the data however raises the privacy concern. The concern particularly affects the display ad market (rather than the search ad market) as the data collection may not align with consumer intention (Budak et al., 2016). Previous papers have examined the economic impact of stricter privacy regulation on the digital display ad market. Goldfarb and Tucker (2011) use large, campaign-level survey data to evaluate the economic impact of the 2002 European Union “Privacy and Electronic Communications Directive” that limited the cross-website tracking. Their calculation indicates that after the stricter privacy regulation, advertisers would have to spend 2.85 times as much advertising to achieve the same level of purchase intent prior to the law. Johnson (2013) focuses on how tracking ban may financially affects the ad market utilizing the real-time bidding technology. His simulated

result indicates that the tracking ban may lead to around 38.5% and 45.5% revenue drops for publishers and advertisers respectively.

In this section I use the difference-in-difference (henceforth DiD) estimate to perform a back-of-the-envelope calculation on the potential economic benefit from the improving frequency management, when the privacy regime switches to allowing perfect cross-website tracking. If cross-website tracking is perfectly allowed, frequency capping can be perfectly implemented across websites – just similar to what advertisers can do in a single medium, e.g. YouTube. In this case, the difficulty of frequency management for a multi-homer would be the same as a single-homer, so there shouldn't be any discount on the valuation of the multi-homer's attention. In other words, multiple websites are just like one super big websites with multiple webpages. In this scenario, there is no multi-homing discount. I evaluate the economic impact by comparing my estimation result from the status quo to the perfect tracking scenario. Since the main multi-homing effect cannot be identified, I cannot assess the price enhancement for all ads. However, I know that the prices of top ads will enhance more (by the DiD estimator) and the question is by how much.

In the previous sub-section I argue that the positive trend that attenuates the DiD estimate comes from advertiser selection – more lasting ads are paired with more viewable ad slots, especially when the multi-homing level is high. The selection may be lessened when the multi-homing level is low as duplication is less of a concern on those websites. In other words, the estimated coefficient of the first-order interaction term ( $-0.782$ ) would be the closest to the unbiased DiD effect. I hence use this figure to perform my calculation. The implicit assumption is that I attribute the full effect of advertiser selection to the second-order interaction term. Note that the figure is likely to serve as the lower bound of the true DiD effect as advertiser selection may still materialize.

In my sample that includes 529 websites, the average multi-homing level (daily and narrowed measure) is 1.88 (see Table 4). Hence, for an average site, if cross-website tracking can be perfectly performed, the log price for a top ad should increase by roughly  $0.782 \times$

$\ln(1.88) = 0.215$  more than a non-top ad, holding other things constant. So even if the prices of non-top ads are not affected by the privacy regulation change (null baseline change), the price of a top ad for an average website can still increase by 21.5% of the non-top ad price. For example, if a top ad sells for \$2 per impression and a non-top ad sells for \$1 per impression in the status quo, the relaxing privacy law would imply that at least the top ad price would reach  $\$2 + \$1 \times 0.215 = \$2.215$ .

Currently in the United States, the annual revenue on the digital display ad market is about \$35 billion. My data comes from BuySellAds that performs the automated transaction service, so my estimation should be most relevant to such segment, which accounts for 40% of the revenue<sup>19</sup>. In my sample, the top ads have an average cpm price of \$6.5 and the non-top ads have an average cpm price of \$3.9. Assume that the ad slot supply of the top ads is roughly equal to the non-top ads, the ad revenue for the top ads is then  $\$35 \times 40\% \times \$6.5 / (\$6.5 + \$3.9) = \$8.75$  billion and the revenue for the non-top ads would be \$5.25 billion. When the privacy regime relaxes, it is then expected that the revenue for top ads may increase at least by  $\$5.25 \times 0.215 = \$1.13$  billion only through the improvement on frequency management, without considering the potential revenue increase in general and holding everything else constant. That is, the \$14 ( $\$35 \times 40\%$ ) billion market of automated transaction service on display ads may otherwise increase its size at least by \$1.13 billion – about 7%, if cross-website tracking is perfectly allowed. If the result can apply to the entire display ad market, the magnitude can be as large as  $\$35 \times 7\% = \$2.45$  billion.

As argued in Budak et al. (2016), strengthening the privacy regulation can simply incentivize media to look for alternative methods of monetization, such as user donation or subscription. An opposite effect works here – relaxing the privacy law may further expand the online display ad market, and thus the “holding everything else constant” essentially puts the calculation to the lower bound of the potential economic impact.

---

<sup>19</sup>The other two segments, direct transaction and real-time bidding, account for 27% and 33% respectively.



## 7 Auxiliary Models

A major concern for adapting the publisher fixed effect model is that it may eliminate too much “good” information, relative to the endogenous or noisy variation. That is, the remedy can be too strong (although conservative) when prices and multi-homing levels are relatively time-invariant. In this section, I propose two other models that, instead of eliminating all publisher time-invariant information, remove only the endogenous error based on careful arguments. The first model utilizes the two-sided market features to develop valid controls. The second model employs several potential instrumental variables for the endogenous multi-homing level. The goal is to refine the comparison between the already balanced treated and control groups and block the endogenous effect from publisher strategic behaviors. Both models indicate that the pattern of  $\Delta V_H$  found in the main model is robust as shown in Table 5.

### 7.1 Two-sided Market OLS

Online readership is fundamentally different from traditional newspapers industry (Gentzkow et al., 2014) and print magazine industry (Shi, 2015). Quantitatively, the scale of consumer multi-homing is much larger and the audience overlapping structure is more complex. Qualitatively, a consumer often allocates time across the purely ad-financed publishers, instead of making subscription decisions. A different choice model is hence developed as follows.

I assume that consumer browsing decision concerning the time allocation across publishers depends on the publisher content, the ad level, and his/her time constraint. The publisher content quality in my context includes both content depth and breadth as both can affect the visiting choices. The ad level means ad supply for an average webpage – higher number of ads can be disturbing or attractive depending on the consumer ad preference. Time allocation across publishers depends on consumer time constraint, which is another idiosyncratic factor. The time spent on websites is directly observed in my data; the content quality and ad level

Table 5: Results for Auxiliary Models

		OLS		IV-2SLS	IV-LIML
	$\ln(cpm)$	$\ln(cpm)$	$\ln(cpm)$	$\ln(cpm)$	$\ln(cpm)$
$T$	0.438*** (4.70)	0.329*** (2.86)	0.300** (2.55)	0.697*** (3.82)	0.732*** (3.66)
$\ln(H)$	1.154*** (3.62)	1.508*** (3.01)	2.202*** (4.05)	3.757*** (3.11)	4.213*** (2.84)
$\ln(H)^2$	-0.568*** (-3.71)	-0.687*** (-2.90)	-0.892*** (-3.27)	-1.591** (-2.54)	-1.814** (-2.36)
$T \times \ln(H)$	-0.396** (-1.99)	-0.429 (-1.44)	-0.912*** (-2.76)	-1.574** (-2.06)	-1.820** (-1.99)
$T \times \ln(H)^2$	0.379*** (3.31)	0.417** (2.10)	0.757*** (3.13)	1.104* (1.92)	1.294* (1.85)
priority	-0.0309 (-0.20)	-0.0520 (-0.29)	-0.0754 (-0.43)		
Ad size	Yes	Yes	Yes	Yes	Yes
Demographics	Yes	Yes	Yes	Yes	Yes
Month FE	Yes	Yes	Yes	Yes	Yes
Genre FE	Yes	Yes	Yes	Yes	Yes
Genre FE $\times$ quality	Yes	Yes	Yes	No	No
Genre FE $\times$ duration	Yes	Yes	Yes	No	No
Cragg-Donald Wald F Value	.	.	.	11.29	11.29
Hansen's J Statistic	.	.	.	0.10	0.13
User base	$\geq 3$	$\geq 10$	$\geq 20$	$\geq 3$	$\geq 3$
Publishers	300	172	116	300	300
Observations	1585	986	714	1585	1585

$t$  statistics in parentheses

\*  $p < 0.10$ , \*\*  $p < 0.05$ , \*\*\*  $p < 0.01$

standard error clustered in the publisher level

are however not. This poses the main challenge to identification, as both can have direct effects on ad prices and the reader composition (Ambrus et al., 2014).

I address the identification challenge by utilizing the rich consumer side information. The idea is to use the observed market share data (in terms of unique visitors) and the average time spent on a webpage to recover the relevant information on the content quality and the ad level. To do this and to reflect the online features, I describe the consumer decision process of information consumption in two stages. In the first stage, a consumer chooses the media publishers to visit based on the content quality and his/her idiosyncratic preference on content. In the second stage, a consumer allocates the time spent on an average webpage of a publisher, based on content quality, ad level, and his/her idiosyncratic preference on ad and time. I then impose critical assumptions to these idiosyncratic parts in both stages, which will be the focus in the coming two sub-sections. The goal is to derive a tractable reduced-form model that accounts for endogeneity.

### ***Advertising Market***

I regard the ad slots sold by different website publishers as heterogeneous products, with the following argument. The ad slots are not perfect substitutes because they may reach different readers (i.e. they are not a homogeneous product), and yet the ad slots are substitutes as they may reach the same reader given that some readers multi-home. This feature distinguishes my model from the “single-homing consumer” model attempting to identify the cross-sided network effect (e.g. Rysman (2004)). That is, there is no “competitive bottleneck” (Armstrong, 2006) as publishers compete in both advertising market and consumer market.

Recall that in my data the publishers sell their ad inventories in an up-front price through the online platform. This means that publishers in my sample do not price discriminate advertisers. The demand system in the advertising market is hence similar to the standard heterogeneous product market (e.g. car), except that advertisers can also patronize any number of publishers, and they will, as long as their valuation is higher than the ad price.

However, given that most consumers do multi-home, when advertisers patronize multiple publishers, their valuation for each impression may fall due to the potential duplicated exposures. For advertisers, this is analogous to buying a second unit of a heterogeneous product, when some functionality has already been fulfilled by the first purchase.

To describe the ad demand, I adapt a linear system. There are two reasons for this. First, as I don't observe advertisers' choices, any demand model (e.g. logit system) requiring knowing information on the choices made is not possible to estimate. However, it is possible to infer the ad level of each publisher from the consumer side data. Second, as my primary goal is not identifying the complete demand structure, the curse of dimensionality (Nevo, 2000) does not materialize. I assume that competitors' ad prices do not correlate with the multi-homing level, once the own ad level is controlled. This means that leaving the price vector of competitors to the residual term in the demand system does not affect the causal interpretation of the multi-homing measures, as demand shifters.

I hence have the following linear specification of publisher  $i$  and ad slot  $j$  for a given time:

$$P_{ij} = \beta_0 + \gamma T_{ij} + f^n(H_i|\delta) + f^n(T_{ij} \times H_i|\lambda) + S'_{ij}\beta_1 + X'_i\beta_2 + \boxed{\pi^a a_i + \pi^q q_i} + \epsilon_{ij}$$

I drop  $t$  for the simplicity on illustrating the model. All notations remain the same as the previous section, except for two new publisher-specific variables. These are two unobserved but endogenous variables:  $a_i$  is the advertising quantity of the webpage, and  $q_i$  is the mean content quality for publisher  $i$ . They are endogenous because the combination of them defines reader's perceived quality on the publisher and thus affects the reader composition, though the effect direction on the multi-homing level is not clear (Ambrus et al., 2014). The idea is that, instead of eliminating every publisher fixed effect, I attempt to measure the two unobserved but endogenous sources and control for them. The assumption for identifying  $\delta$  is that by removing the two endogenous sources the rest shock  $\epsilon_{ij}$  becomes exogenous.

Such assumption may be too strong given the endogeneity nature of market prices. The assumption for identifying  $\lambda$  is however less strict. It only requires that the bias works in the same manner to the treated and to the control.

### ***Consumer Market***

To infer  $a_i$  and  $q_i$ , I assume that a consumer consumes media in two stages. In the first stage, a consumer chooses to patronize some online media publishers based on the content quality and his/her idiosyncratic preference. The idea is to use the consumer market demand to infer the publisher mean quality. Given that it is common for consumer to patronize multiple publishers, potentially I may follow Gentzkow (2007) to construct  $2^J$  “bundles,” where each bundle can contain  $0, 1, 2, \dots, J$  publishers, and estimate a choice model of those bundles. The method is general to both substitution and complementarity patterns, but the issue in the online media is that  $J$  can be large. The smallest genre in my matched data sample contains 13 publishers and can be as large as 2249. The scale problem makes the estimation of mean quality computationally non-executable.

Another much simpler approach is to assume that consumers treat each publisher as a completely different product. Although this may look arbitrary at the first glance, I argue that this is reasonable for describing the online readership. A fundamental difference of online media, unlike a magazine or a newspaper, is that the information consumption can break down into separable webpages. Even visiting one webpage would still count as a visitor, which aggregately constitutes the consumer demand for the publisher. This means that for two publishers, if there is a least one webpage that provides completely independent information, the two can be regarded as different products. In my data, I observe the monthly number of unique visitor of each publisher. It should be fair to say that the at least parts of the webpages of any two publishers over a month should convey independent information.

In my empirical model, I adapt a simple logit model and treat each publisher’s content provision (over a month) separately as different products. Hence, consumers make separate binary choice for each publisher. The mean quality can then be inferred using the “Berry

inversion” in the simplest form as the following.

$$q_i = \log\left(\frac{n_i}{m_g - n_i}\right)$$

where  $n_i$  is the number of unique visitors of publisher  $i$  and  $m_g$  is the number of unique visitors of genre  $g$  (the market size). I use the sample in comScore to calculate the measure. Essentially, the content quality is approximated by a function of the number of unique visitors. Intuitively, the higher number of unique visitors comes from either the direct visits from loyal customers, from the organic search results, or the advertised search results, and it should reflect the content quality of a publisher.

In the second stage, a consumer decides how much time to spend on the selected media publishers, based on content quality, ad level, time constraint, and his/her idiosyncratic preference. Essentially, I assume that the ad level can only affect the “depth” when consumers determine the time spent on a publisher, but not the “breadth” when consumers choose which publishers to visit. The assumption is not uncommon in literature and it is often referred to consumer’s “passive” beliefs on the ad level (Gabszewicz and Wauthy, 2005; Anderson et al., 2015; Anderson and Jullien, 2015). As consumers do not observe the ad level before visiting the publisher, they form an expectation to it. The belief system is passive if consumers’ expected ad level does not vary as publishers chooses some different levels. The result of this assumption is that consumer’s decision on patronizing publishers does not depend on the actual ad level. In the online context, when consumers are costless to access and leave webpages, the role of this expectation on choosing media outlets can be further minimized.

I now describe the model specification of the second stage. I define ad level  $a_i$  as the number of ads in a webpage. The ad level may negatively, neutrally, or positively affect consumer time allocation on that webpage, depending on the preference genre-specific parameter  $\alpha_g$ , followed Kaiser and Song (2009). I denote  $duration_{ki}$  to consumer  $k$ ’s optimal

time allocation to an average webpage of publisher  $i$ , and it can be presented by a linear relationship of the publisher content main quality  $q_i$ , the ad level  $a_i$ , the genre-specific time factor  $\omega_g$ , and the idiosyncratic preferences  $\alpha_{kg}$  and  $\omega_{kg}$  as follows:

$$duration_{ki} = \phi_g q_i + (\alpha_g + \alpha_{kg}) a_i + (\omega_g + \omega_{kg})$$

The intuition is that different genres may have different features. For some, consumers are more sensitive (in terms of time spent) to quality ( $\phi_g$ ). For some, consumers may like or dislike ads more ( $\alpha_g$ ). The parts that are unexplained by the quality and ad level is captured in the genre-specific time factor ( $\omega_g$ ). Within each genre, consumer may preserve different patterns measured by  $\alpha_{kg}$  and  $\omega_{kg}$ . Those idiosyncratic parts may become immaterial when the duration measure is aggregated to the publisher level, under the following conditions, where  $\zeta_k$  denotes the logistic idiosyncratic random vector in the first stage.

$$E(\alpha_{kg}|\zeta_k) = E(\alpha_{kg}) = 0 \quad \text{and} \quad E(\omega_{kg}|\zeta_k) = E(\omega_{kg}) = 0$$

The conditions mean that within a genre, the consumer choices on publishers do not correlate to idiosyncratic ad preference ( $\alpha_{kg}$ ) structure and time constraint factor ( $\omega_{kg}$ ). That is, a publisher is unlikely to gather a group of visitors who, on average, like (dislike) ads more or have more leisure time, than the other publishers' visitors in the same genre market<sup>20</sup>. A potential augmentation is to utilize the detailed demographic data and include them into the linear specification of  $duration_{ki}$ . But doing so does not change the final specification of my reduced-form model. As my purpose is not to identifying all those demographic coefficients, I would rather keep the linear specification parsimonious. The approximation of the ad level

---

<sup>20</sup> As the news media example illustrated in White and Weyl (2016), highbrow and lowbrow publishers may attract different readers who have different ad tastes. In my case the distinction of highbrow and lowbrow is minimum as all my publishers are purely ad-financed (no subscription fee).

of publisher  $i$  in genre market  $g$  thus can be derived:

$$duration_i = E(duration_{ki}|\zeta_k) = \phi_g q_i + \alpha_g a_i + \omega_g$$

Rearrange:

$$a_i = \frac{1}{\alpha_g}(duration_i - \phi_g q_i - \omega_g)$$

Plug the equation to the econometric specification, and the regression model becomes:

$$P_{ij} = \beta_0 + \gamma T_{ij} + f^n(H_i|\delta) + f^n(T_{ij} \times H_i|\lambda) + S'_{ij}\beta_1 + X'_i\beta_2 + \boxed{\frac{\pi^a}{\alpha_g}duration_i + (\pi^q - \frac{\phi_g \pi^a}{\alpha_g})q_i - \frac{\pi^a \omega_g}{\alpha_g}} + \epsilon_{ij}$$

Adapting the publisher-month panel structure and rewrite the model to the following reduced-form specification:

$$P_{ijt} = \beta_0 + \gamma T_{ijt} + f^n(H_{it}|\delta) + f^n(T_{ijt} \times H_{it}|\lambda) + S'_{ijt}\beta_1 + X'_{it}\beta_2 + \boxed{\pi_g^a duration_{it} + \pi_g^q q_{it} + \rho_g} + \epsilon_{ijt}$$

The specification essentially says that the unobserved but endogenous aggregate variation from the ad level  $a_{it}$  and the content quality  $q_{it}$  can be captured by the average webpage duration, the market share, the genre interaction effects on the two variables, and the genre fixed effect<sup>21</sup>. The results are shown in the first three columns in Table 5. Compare to the main model using the publisher fixed effect, the pattern of the estimation results is similar.

### ***Local Greater-reach “priority”***

Finally, I explain the “priority” measure I include in the OLS regression. The measure of “priority” is used to capture the greater-reach premium in the complex composition of website visitors. As in the simple two-publisher model, (Athey et al., 2016) show that the publisher with larger user base (unique visitors) can receive the priority from advertisers

---

<sup>21</sup>I assume that  $\pi^a$ ,  $\omega_g$ , and  $\alpha_g$  are time invariant.



for showing their ads. In their advertiser mix-homing equilibrium, this prioritized publisher receives the full ad demand, but the other one only captures the partial ad demand as some advertisers cannot tolerate the efficiency from duplicated impressions. Hence, some advertisers patronize both publishers while some only patronize the larger one – to reach as many audience as possible without losing efficiency.

In the online genre market, however, there are extensive amounts of publishers, potentially without a single dominator in terms of reach. Advertisers are likely to choose a set of publishers, instead of one or all, to show their ads. The measure of “priority” is meant to capture this greater-reach premium, in addition to the number of unique visitors that applies most relevantly in the simple two-publisher situation. The logic is as the following. To optimize the reach without too much wasteful duplication, in a given genre market advertisers can break the websites into groups based on their audience overlapping structure. The ideal division is to have several groups that do not overlap with others. Anything close to that can be an effective strategy, similar to the “key players” problem in social networks (An and Liu, 2016). For each overlapping group, advertisers then pick one most extensive website to show their ad. In this sense, in each local group there is one website that can receive this premium, even if this local group has smaller audience base than websites in other groups. Therefore, rather than an absolute measure on reach, I develop a relative measure to account for the additional complexity rooted from the scale of Internet.

I define the “neighbors” of website  $i$  as those websites who share audiences with  $i$ . By definition,  $i$  may capture any percentage (greater than 0) of its neighbor’s audience. The higher the percentage, the more dominant  $j$  is to its neighbor. The “priority” measure is the aggregate level of this dominance. I name so as more dominance can translate to the priority of ad demand, which can offset the multi-homing discount. To an extreme case, a website may 100% dominate all its, say five, neighbors. The “priority” measure is then five, meaning that it covers the web traffics of all its five neighbors.

Despite that each of them represents a specific channel, the “multi-homing” and “priority”

measures can be correlated and they do in my data. In the two-publisher case illustrated in (Athey et al., 2016), the two measures are perfectly aligned. The publisher with greater reach receives the priority and the same publisher must have lower multi-homing level. In the reality where there are many more publishers in the same genre market, how the two measures are related is however an empirical question. My sample indicates that, in fact, many locally dominant websites (high “priority”) also have high multi-homing level (the correlation coefficient is 0.42), meaning that they may escape from the undesired multi-homing discount by possessing greater local reachability. When including the measure into the previous models, however, it does not alter the pattern of multi-homing effect.

## 7.2 Instrumental Variables

Instead of retrieving valid measures for publisher strategic variables, alternatively I may search for instrumental variables (IVs) that are exogenous to them. I consider two IVs for the multi-homing levels. First, I use the consumer multi-homing behavior on the websites that are far from the examining genre market. The implicit assumption is that there is a general tendency for consumer to multi-home (the relevance condition), but the unobserved publisher strategies do not affect the consumer multi-homing behavior out of the genre market (the exclusion condition). Specifically, The suspected endogenous variable is the multi-homing level within the narrow genre ( $H^{narrow}$ ), and the instrument is the multi-homing level outside the broad genre market ( $H^{out} = H^{universal} - H^{broad}$ ). The second IV is user total time spent online. I assume that such time spent may be correlated to the consumer multi-homing behavior, but is free from the strategy of a single publisher.

In the estimation, there are four endogenous variables ( $\ln(H)$ ,  $\ln(H)^2$ ,  $T \times \ln(H)$ ,  $T \times \ln(H)^2$ ). I use the two IVs to create the corresponding functional forms to each. I also augment the set of IVs by using the predict  $\ln(H)$  in the first stage to create one more square term and one more interaction term. That is, I use the ten exogenous variables to instrument for the four endogenous variables.

As shown in Table 5 the 4th and 5th column, the estimated coefficients in the interaction terms have a similar pattern to other models. Both the 2SLS and LIML procedure report similar results, indicating that the weak identification problem is mild. The first stage F value is over ten and they pass the overidentification test. Note that, however, this is not true when restricting the user base thickness to a higher level. The estimation results for user base  $\geq 10$  and  $\geq 20$  are very different between the 2SLS and LIML procedures and the first stage F value for these cases are around 4.

## 8 Conclusion

In this paper, I empirically examine whether consumer multi-homing behavior leads to over impressions of online display ads, and whether the duplicated impressions can decrease advertisers' valuation to ad slots. Using a difference-in-difference strategy based on the ad locations in a webpage, I find that consumer multi-homing behavior within a narrow genre has a negative impact on the ad prices in that genre market. The result is robust to sample selection and different reduced-form models. With careful and plausible assumptions, I interpret that the finding implies that consumer multi-homing indeed causes undesired duplication that lowers advertisers willingness-to-pay.

The paper suggests that advertisers still face a challenge on frequency management under the current self-regulatory privacy principles. The current privacy principles allow frequency tracking within a website as well as within an ad network. As ad networks possess an advantage on cross-website tracking from the current principles, a further expanding on the networks or any merger and acquisition among ad networks may increase advertisers' valuation on an average ad impression. When the "third-party" tracking is perfectly achieved or allowed, my estimation results suggest that the ad revenue can increase by at least 7%, which translates to \$1.13 to \$2.45 billion annually, depending on the the definition of the applicable market (automated transaction segment only or the whole online display ad market).

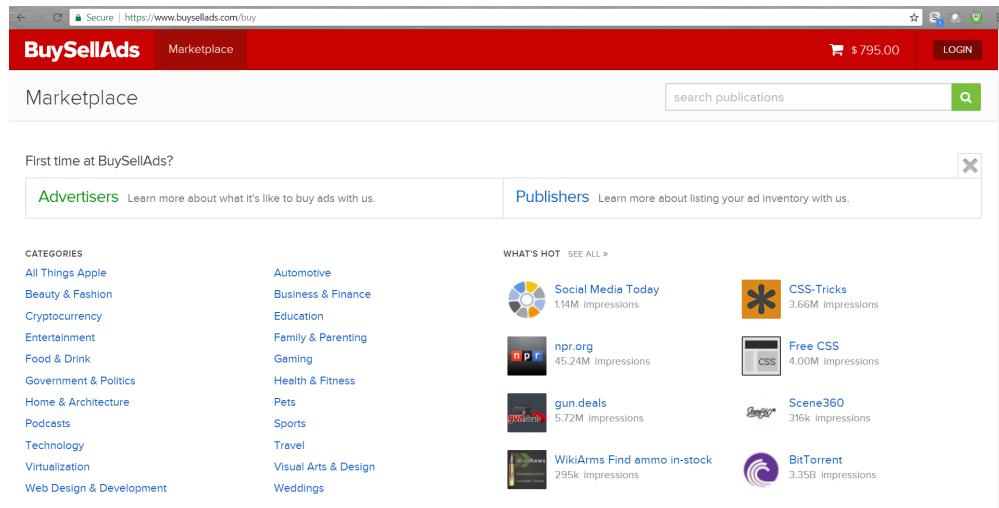
The paper also suggests that the existing economic analyses based on the multi-homing-duplication channel can be applied to the digital display ad market. When the privacy protection becomes stricter, the impact of consumer multi-homing can be more significant: online publishers may tend to either produce niche contents to attract less multi-homing readers or produce wide contents to attract less-multi-homing advertisers.

## References

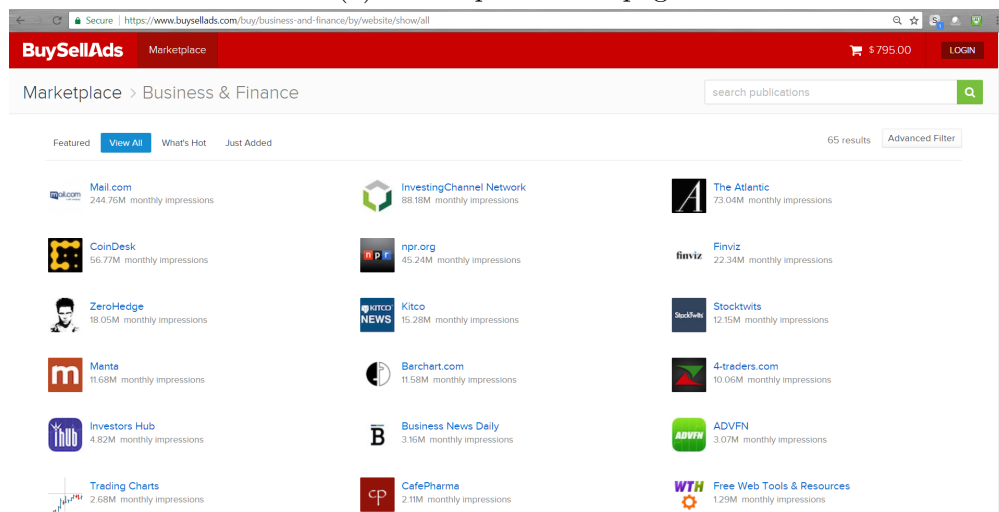
- Ambrus, A., Calvano, E., and Reisinger, M. (2014). Either or both competition: a “two-sided” theory of advertising with overlapping viewerships. *Working paper*.
- An, W. and Liu, Y.-H. (2016). keyplayer: An r package for locating key players in social networks. *R Journal*, 8(1):257–268.
- Anderson, S. P. and Coate, S. (2015). Market provision of broadcasting: A welfare analysis. *Review of Economic Studies*, 72(4):947–972.
- Anderson, S. P., Foros, Ø., and Kind, H. J. (2015). Competition for advertisers and for viewers in media markets. *Working paper*.
- Anderson, S. P. and Jullien, B. (2015). The advertising-financed business model in two-sided media markets. *Working paper*.
- Armstrong, M. (2006). Competition in two-sided markets. *RAND Journal of Economics*, 37(3):668–691.
- Athey, S., Calvano, E., and Gans, J. S. (2016). The impact of consumer multi-homing on advertising markets and media competition. *Working paper*.
- Budak, C., Goel, S., Rao, J., and Zervas, G. (2016). Understanding emerging threats to online advertising. *Working paper*.
- Chandra, A. (2009). Targeted advertising: The role of subscriber characteristics in media markets. *Journal of Industrial Economics*, 57(1):58–84.
- Chandra, A. and Kaiser, U. (2014). Targeted advertising in magazine markets and the advent of the internet. *Management Science*, 60(7):1829–1843.
- Flexman, S., Goel, S., and Rao, J. M. (2016). Filter bubble, echo chambers, and online news consumption. *Public Opinion Quarterly*, 80:298–320.
- FTC (2009). Federal trade commission staff report: Self-regulatory principles for online behavioral advertising. [link](#).
- Gabszewicz, J. J. and Wauthy, X. Y. (2005). Two-sided markets and price competition with multi-homing. *Working paper*.

- Gentzkow, M. (2007). Valuing new goods in a model with complementarity: Online newspapers. *American Economic Review*, 97(3):713–744.
- Gentzkow, M. and Shapiro, J. M. (2011). Ideological segregation online and offline. *Quarterly Journal of Economics*, 126(4):1799–1839.
- Gentzkow, M., Shapiro, J. M., and Sinkinson, M. (2014). Competition and ideological diversity: Historical evidence from us newspapers. *American Economic Review*, 104(10):3073–3114.
- Goldfarb, A. and Tucker, C. (2011). Privacy regulation and online advertising. *Management Science*, 57(1):57–71.
- Interactive Advertising Bureau (2013). Programmatic and automation the publishers perspective. link. accessed September 2017.
- Jeziorski, P. (2014). Effects of mergers in two-sided markets: The us radio industry. *American Economic Journal: Microeconomics*, 6(4):35–73.
- Johnson, G. A. (2013). The impact of privacy policy on the auction market for online display advertising. *Working paper*.
- Kaiser, U. and Song, M. (2009). Do media consumers really dislike advertising? an empirical assessment of the role of advertising in print media markets. *International Journal of Industrial Organization*, 27(2):292–301.
- Nevo, A. (2000). A practitioner’s guide to estimation of random-coefficients logit models of demand. *Journal of Economics and Management Strategy*, 9(4):513–548.
- Rysman, M. (2004). Competition between networks: A study of the market for yellow pages. *Review of Economic Studies*, 71(2):483–512.
- Shi, C. M. (2015). Catching (exclusive) eyeballs: Multi-homing and platform competition in the magazine industry. *Working paper*.
- Steiner, P. O. (1952). Program patterns and the workability of competition in radio broadcasting. *Quarterly Journal of Economics*, 66(2):194–223.
- Sthmeier, T. and Wenzel, T. (2012). Regulating advertising in the presence of public service broadcasting. *Review of Network Industries*, 111(2):1–21.
- White, A. and Weyl, G. (2016). Insulated platform competition. *Working Paper*.

# Figures




(a) Marketplace Homepage



(b) List of Publishers in Business & Finance

Figure 2: BuySellAds (BSA) Marketplace






The Atlantic  
The Atlantic Monthly Online | [theatlantic.com](http://theatlantic.com)

Channel: [Government & Politics & Business & Finance](#)  
Member Since October 2014

71,300,000  
Monthly Impressions

1,596,089  
Followers

The Atlantic Monthly's home on the Internet, featuring current issues online alongside web-specific content on travel, literature, politics, and digital culture.

Ask a question

(a) Ad Inventories of the Atlantic

Website (CPM)

Targeting
Scheduling

Category Targeting

Search...

Education -\$10.00  
Entertainment -\$4.00  
Health -\$4.00  
International -\$4.00  
National -\$10.00  
Politics -\$4.00  
Technology -\$10.00

Selected Criteria

CATEGORY TARGETS

Business -\$10.00

GEO TARGETS

United States (Country) -\$3.00

**Top Right**  
300 x 600 Top Right

932,000  
Est Impressions

✓  
Available

**\$25.00 (+\$13.00)**  
CPM

+

Add To Cart

**Top Right**  
300 x 250 Top Right

883,000  
Est Impressions

✓  
Available

**\$17.85 (+\$13.00)**  
CPM

+

Add To Cart

**Leaderboard**  
728 x 90 Top Center  
CURRENTLY SOLD OUT

1,000  
Est Impressions

Unavailable

\$17.00 (+\$13.00)  
CPM

Waiting List

**In Content**  
300 x 250 Center

4,281,000  
Est Impressions

✓  
Available

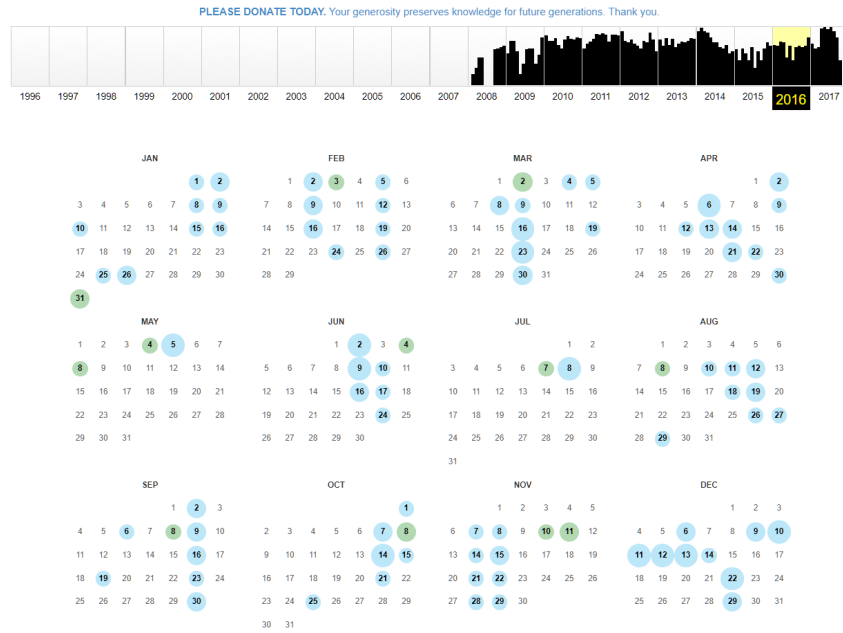
**\$16.00 (+\$13.00)**  
CPM

+

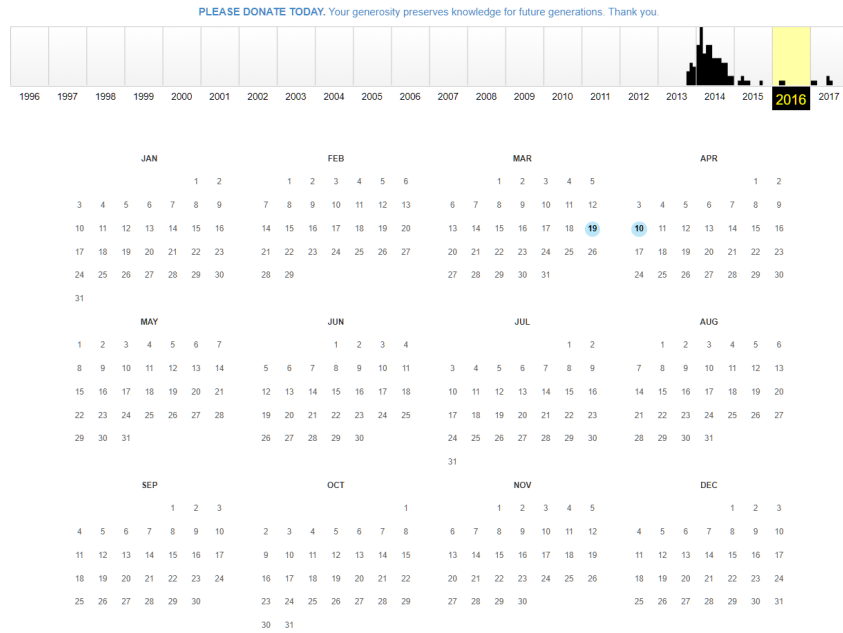
Add To Cart

(b) Targeting Options of the Atlantic

Figure 3: BuySellAds (BSA) User Interface



(a) BSA Marketplace Pages Archived



(b) BSA Business & Finance Pages Archived

Figure 4: Webpages Archived by Dates (Blue: webpage preserved; Green: webpage redirected)