



**HAL**  
open science

# **L1 chimeric transcripts are expressed in healthy brain and their deregulation in glioma follows that of their host locus**

Marie-Elisa Pinson, Franck Court, Aymeric Masson, Yoan Renaud, Allison Fantini, Ophélie Bacoœur-Ouzillou, Marie Barriere, Bruno Pereira, Pierre-Olivier Guichet, Emmanuel Chautard, et al.

► **To cite this version:**

Marie-Elisa Pinson, Franck Court, Aymeric Masson, Yoan Renaud, Allison Fantini, et al.. L1 chimeric transcripts are expressed in healthy brain and their deregulation in glioma follows that of their host locus. *Human Molecular Genetics*, 2022, 31 (15), pp.2606-2622. 10.1093/hmg/ddac056 . hal-03631345

**HAL Id: hal-03631345**

**<https://hal.science/hal-03631345>**

Submitted on 5 Apr 2022

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

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

1 **L1 chimeric transcripts are expressed in healthy brain and their deregulation in glioma follows**  
2 **that of their host locus.**

3

4 **Marie-Elisa Pinson <sup>1§</sup>, Franck Court <sup>1§</sup>, Aymeric Masson <sup>1</sup>, Yoan Renaud <sup>1</sup>, Allison Fantini <sup>1</sup>,**  
5 **Ophélie Bacoer-Ouzillou <sup>1</sup>, Marie Barriere <sup>1</sup>, Bruno Pereira <sup>2</sup>, Pierre-Olivier Guichet <sup>3</sup>,**  
6 **Emmanuel Chautard <sup>4,5</sup>, Lucie Karayan-Tapon <sup>3,6,7</sup>, Pierre Verrelle <sup>8,9,10</sup>, Philippe Arnaud <sup>1</sup>,**  
7 **Catherine Vaurs-Barrière <sup>1\*</sup>.**

8 1 Université Clermont Auvergne, CNRS, Inserm, iGReD, F-63000 Clermont-Ferrand, France.

9 § co-first authors: Both authors contribute equally to this work

10 2 Biostatistics Department, Délégation à la Recherche Clinique et à l'Innovation, Clermont-Ferrand  
11 Hospital, Clermont-Ferrand 63003, France

12 3 CHU de Poitiers, Cancer Biology Department, Poitiers 86021, France

13 4 Université Clermont Auvergne, INSERM, U1240 IMoST, Clermont-Ferrand 63011, France

14 5 Centre Jean PERRIN, Pathology Department, Clermont-Ferrand 63011, France

15 6 INSERM, U1084, Poitiers 86021, France

16 7 Université de Poitiers, Poitiers 86000, France

17 8 INSERM, U1196 CNRS UMR9187, Curie Institute, Orsay 91405, France

18 9 Radiotherapy Department Curie Institute, Paris 75005, France

19 10 Université Clermont Auvergne, Clermont-Ferrand 63000, France

20

21 \*Corresponding author:

22 Vaurs-Barrière Catherine

23 iGReD

- 1 28 Place Henri Dunant
- 2 63000 Clermont-Ferrand FRANCE
- 3 phone: +33 (0)4 73 17 81 78
- 4 Fax : +33 (0)4 73017 83 83
- 5 email: [catherine.barriere@uca.fr](mailto:catherine.barriere@uca.fr)
- 6

1 **Abstract** (250 words)

2 Besides the consequences of retrotransposition, long interspersed element 1 (L1) retrotransposons can  
3 affect the host genome through their antisense promoter. In addition to the sense promoter, the  
4 evolutionarily recent L1 retrotransposons, which are present in several thousand copies, also possess  
5 an anti-sense promoter that can produce L1 chimeric transcripts (LCT) composed of the L1 5' UTR  
6 followed by the adjacent genomic sequence. The full extent to which LCT expression occurs in a  
7 given tissue and whether disruption of the defense mechanisms that normally control L1  
8 retrotransposons affects their expression and function in cancer cells, remain to be established.

9 By using CLIFinder, a dedicated bioinformatics tool, we found that LCT expression was widespread  
10 in normal brain and aggressive glioma samples, and that approximately 17% of recent L1  
11 retrotransposons, from the L1PA1 to L1PA7 subfamilies, were involved in their production.  
12 Importantly, the transcriptional activities of the L1 antisense promoters and of their host loci were  
13 coupled. Accordingly, we detected LCT-producing L1 retrotransposons mainly in transcriptionally  
14 active genes and genomic loci. Moreover, changes in the host genomic locus expression level in  
15 glioma were associated with a similar change in LCT expression level, regardless of the L1 promoter  
16 methylation status.

17 Our findings support a model in which the host genomic locus transcriptional activity is the main  
18 driving force of LCT expression. We hypothesize that this model is more applicable when host gene  
19 and LCT are transcribed from the same strand.

20

21

## 1 **Introduction**

2 Long interspersed nuclear element 1 (LINE-1 or L1) is a class of autonomous transposable elements  
3 that is mobilized *via* an RNA intermediate by a copy-and-paste mechanism named retrotransposition.  
4 L1 retrotransposons represent approximately 17% of the human genome (1), and 516,000 copies have  
5 been annotated in the human reference genome. Seven thousands of them are full length (~6kb) L1  
6 retrotransposons among which only ~100 can still retrotranspose *in vivo* and belong to the most recent  
7 evolutionary L1 subfamily, specific to the human species (L1Hs or L1PA1 elements) (2).

8 Besides the consequences of retrotransposition (*e.g.* insertional mutagenesis, creation of new  
9 alternative splicing sites, promoting sequence transduction from the donor to new insertion sites) (3),  
10 L1 retrotransposons can also interfere with the transcriptional activity of the surrounding genomic  
11 sequences from their bidirectional promoter. Specifically, in addition to the sense promoter, the  
12 5'UTR of evolutionarily recent L1s also contain an anti-sense promoter (L1-ASP) that can produce  
13 transcripts from the 5' UTR in antisense orientation to the adjacent genomic region. These transcripts  
14 are called L1 chimeric transcripts (LCTs). The L1-ASP was initially described in the L1PA1  
15 subfamily (4), but sequence homology and functional analyses suggest that the L1-ASP appeared first  
16 in the L1PA6 subfamily and is active in the more recent subfamilies (L1PA1 to L1PA6) (5,6). L1-ASP  
17 activity has been detected also in older L1 families, up to L1PA8, despite the lower sequence  
18 homology (6,7). Moreover, it has been demonstrated that the L1 5'UTR transcribed in the antisense  
19 orientation from the L1-ASP contains an open reading frame termed ORF0 (8). Its translation gives  
20 rise to a short peptide, or more rarely, to fusion proteins with proximal exons (8). This suggests that at  
21 least some LCTs might encode proteins.

22 In many cancer types, disruption of the defense mechanisms that normally control L1 retrotransposons  
23 promotes their expression and mobilization (9,10). Specifically, cancer-associated aberrant DNA  
24 hypomethylation of the L1 5'UTR correlates with L1 transcriptional activation (10,11). Therefore, it  
25 has been proposed that L1 retrotransposons can contribute to oncogenesis by somatic neo-  
26 transposition (that can be a driver event) (11,12), and also by aberrant LCT transcription, a  
27 phenomenon called “onco-exaptation” (13). For instance the *L1-MET* LCT initiates in the second  
28 intron of the gene encoding the tyrosine protein kinase MET (*c-MET*) (14). In various cancer types,

1 *L1-MET* LCT expression is inversely correlated with methylation of its L1 5'UTR, and might have  
2 oncogenic functions (15–18). Similarly, it has been suggested that LCT13, a tumor-specific long non-  
3 coding LCT, induces repression of the 300kb distant tumor suppressor gene *tissue factor pathway*  
4 *inhibitor 2 (TFPI-2)* in various malignancies (19).

5 These observations stress that recent L1s (L1PA1 to PA8) with a full length 5'UTR represent an  
6 abundant source of alternative promoters that can disrupt the expression of nearby genes in cancer.  
7 This highlight the need to identify LCTs in a systematic manner in the different cancer types. Until  
8 now, LCTs have been characterized mainly using bioinformatics approaches. For instance, dbEST has  
9 been queried to identify spliced ESTs in which an L1 5'UTR sequence in antisense orientation is  
10 spliced with a gene exon (4,7,14,17,18). In 2009, Cruickshanks and Tufarelli (20) developed a  
11 dedicated molecular approach (called L1 Chimera Display) that allowed identifying 18 LCTs (intronic  
12 or intergenic) in breast cancer samples and cell lines. Altogether, these studies identified 161 recent L1  
13 retrotransposons implicated in the production of LCTs that are expressed in various normal tissues  
14 and/or cancer types or cell lines, and among which some might be cancer-specific (19,20). However,  
15 these studies had very specific experimental criteria (position of the L1 primer in the LCT and  
16 selection only of spliced exon-containing LCT ESTs) and limitations (*i.e.* absence of weakly  
17 expressed transcripts in EST libraries). Therefore, to gain insights into the LCT expression profile in a  
18 given tissue and into its relevance to cancer development and progression, we and others developed  
19 tools to identify transposable element- derived chimeric transcripts from RNA-seq data. These tools  
20 include the LIONS suite (21), a tool to identify transposable element-derived oncogene transcripts  
21 (22), and Chimeric LINE Finder (CLIFinder), a bioinformatics tool we developed to identify chimeric  
22 reads that correspond to potential LCTs in RNA-seq data (23).

23 Glioma is the most frequent primary malignant brain tumor. The most aggressive forms, mainly  
24 glioblastoma multiform, are molecularly defined by the presence of wild type isocitrate dehydrogenase  
25 1 and 2 (*IDHwt*) (24,25). *IDHwt* glioma is a major source of morbidity and mortality, because it is  
26 almost always fatal. Despite aggressive treatment, these cancers are highly recurrent, leading to a  
27 median survival time after diagnosis that does not exceed 18 months. It has been proposed that

1 therapeutic resistance and tumor relapse rely on a subpopulation of cells within the tumor with stem  
2 cell characteristics, called glioma stem cells (GSC) (26,27).

3 In the present study, we used CLIFinder to explore the LCT landscape in *IDHwt* gliomas and GSC cell  
4 lines to evaluate whether they warrant future investigations as potential players in aggressive glioma  
5 biology. We found that LCTs are expressed in healthy brain and that their transcriptional  
6 activity/deregulation in glioma follows that of their host locus.

7

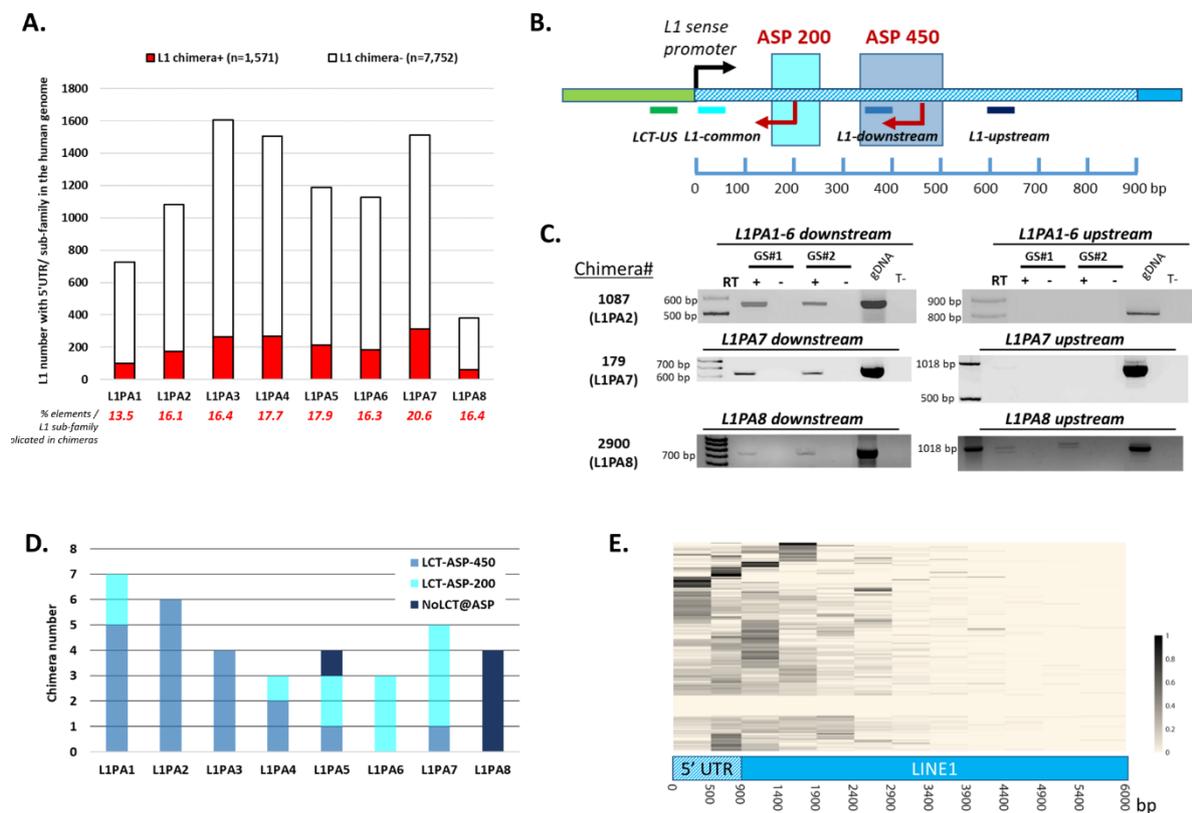
## 8 **Results**

### 9 **Most chimeras including L1PA1 to L1PA7 elements correspond to LCTs**

10 First, we generated Ribo-Zero stranded paired-end RNA-sequencing (RNA-seq) data to analyze the  
11 whole transcriptome (total RNA) of three healthy brain and 8 *IDHwt* glioma samples with CLIFinder.  
12 As the L1-ASP is active in the more recent subfamilies (5,7,19), we used this approach to identify  
13 chimeric transcripts involving one of the 9,123 elements with a 5'UTR from the L1PA1 to L1PA8  
14 subfamilies (Table S1). Although we obtained 90M of reads, each RNA-seq dataset covered only  
15 partially the whole LCT transcriptome of that sample. Accordingly, the number of detected LCTs  
16 increased with the number of samples analyzed, indicating that for extensive genome-wide LCT  
17 identification sequencing data from several samples must be pooled (Figure S1). Finally, by pooling  
18 the RNA-seq data from all 11 samples, we could identify a total of 1,571 chimeras (Table S2), and  
19 observed that all eight subfamilies were involved in their production. The percentage of L1 elements  
20 implicated in chimera production ranged from 13.5% for L1PA1 to 20.6% for L1PA7 (Fig. 1a).

21 LCT transcription initiation from L1-ASP was initially mapped between positions +300 and +500 of  
22 the L1PA1 5'UTR (4), but can also occur between position +150 and +250 (7,20). Chimeras detected  
23 with our approach may initiate from one of these two regions, or may also correspond to larger  
24 transcripts that encompass the L1 5'UTR in antisense orientation. Therefore, we next evaluated the  
25 proportion of chimeras that corresponded to LCTs (*i.e.* they initiated from one of these two L1-ASP  
26 regions). To this aim, we designed a RT-PCR 5' L1 walking approach using three primers, located at  
27 different positions of the L1 5'UTR, and one LCT-specific primer (LCT-Unique Sequence = LCT-US)

1 , located in the unique adjacent sequence, to delineate the area where each chimera initiated (*i.e.* +200  
 2 or +450 in the 5'UTR or upstream) (Fig. 1b).



3

#### 4 **Figure 1: LCTs are produced from the recent L1PA1 to L1PA7 L1 subfamilies**

5

6 We tested 36 chimeras (7 from L1PA1, 6 from L1PA2, 4 from L1PA3, 3 from L1PA4, 4 from L1PA5,  
 7 3 from L1PA6, 5 from L1PA7, and 4 from L1PA8) in two glioma samples (GS#1 and #2). We could  
 8 amplify 19 of them using the two primers located in the first 500bp of the 5'UTR, but not with the  
 9 “upstream” primer, for instance chimera 1087 and 179 (Fig. 1b-d). This pattern was representative of a  
 10 transcription start site (TSS) located in the classical +450 L1-ASP region. For 12 chimeras, we  
 11 obtained an amplification product with the most “downstream” L1 primer (called L1-common), but  
 12 not when using the L1-downstream and L1-upstream primers, flanking the +450 ASP position (Fig.  
 13 1b-d). This suggests that these chimeras might initiate in the +200 L1-ASP region. We could amplify  
 14 the last five chimeras (including the four L1PA8 chimeras tested), such as chimera 2900, using the  
 15 three primers, indicating they do not initiate from the 5'UTR L1-ASP regions, but from an  
 16 undetermined upstream region (NO\_LCTs). This analysis confirmed that the four tested L1PA8  
 17 chimeras corresponded to larger transcripts, and that 31 of the other 32 chimeras (96.8%), which were

1 associated with L1 elements that belonged to the L1PA1 to L1PA7 subfamilies, were LCTs the  
2 transcription of which started at one of the two already described L1-ASP regions.

3 To evaluate whether this observation can apply to a whole population of chimera, we also analyzed by  
4 Ribo-Zero RNA-seq two *IDHwt* glioma stem cell lines (GSC1 and GSC2). By pooling the data of  
5 these two samples, we could identify by CLIFinder 226 chimeras that involved L1 elements from the  
6 L1PA1 to L1PA7 subfamilies. Then, to determine the proportion of chimeras that corresponded to  
7 LCTs, we mapped their TSS along the L1 sequence using a long-read single molecule real-time  
8 (SMRT) sequencing dataset generated from three GSC lines (GSC1, GSC2 and GSC6). Among the  
9 226 chimeras identified by CLIFinder, 22 were not detected in this dataset. The remaining 204 had at  
10 least one TSS within the L1 sequence, including the L1 coding region. Nevertheless, we observed a  
11 TSS hotspot within the 900 bp of the L1 5'UTR (Fig. 1e). Specifically, 98% (201/204) of these  
12 chimeras had one or more 5' ends in the L1 5'UTR and 85% (173/204) within the first 500bp region  
13 that contains the two described ASP regions. This original observation highlights that besides  
14 initiation from the already described L1-ASP region in the 5'UTR, some LCTs can also have  
15 secondary initiation site(s) within the L1 coding sequence, and older L1 elements may be involved in  
16 LCT transcription.

17 These findings suggest that the majority of the 1,509 chimeras implicating these seven subfamilies  
18 (L1PA1 to L1PA7) may correspond to LCTs that initiate at the L1-ASP. For the rest of this study, we  
19 will consider only L1PA1 to L1PA7 as recent L1 subfamilies.

20 Previous studies based on EST querying showed that among the 8,744 recent L1 elements with a  
21 5'UTR sequence from the L1PA1 to L1PA7 subfamilies, 161 (1.9%) are associated with LCT  
22 transcription in various human normal tissue and cancer types and cell lines (4,7,14,17,18,20). By  
23 comparing the genomic positions of the 1,509 LCT-producing L1 elements identified by CLIFinder  
24 and of these 161 "LCT-EST" loci, we determined that 65 L1 elements corresponded to already known  
25 LCT-EST loci (Figure S2) among which 10 were previously described as producing LCT-ESTs in  
26 normal and malignant brain samples ((7); Figure S3).

27 In conclusion, in normal brain and glioma samples, our approach identified 40% of the 161 L1  
28 elements previously associated with LCT-EST production. It also identified 1,444 new recent L1

1 elements implicated in LCT transcription, thus implicating 17.25% of human L1PA1 to L1PA7  
2 retrotransposons (with a 5'UTR) in LCT expression in *IDHwt* glioma and normal brain.

3

#### 4 **Most LCTs are unspliced at their 5'end and can be polyA and non-polyA transcripts**

5 The size of the 1,509 LCTs ranged from 132 bp to 45,574 bp. The median size of the LCTs identified  
6 by only 1 read (n=453) was 306 bp, and 77% of them were smaller than 400 bp. As the L1 antisense  
7 5'UTR contains two donor splicing sites, we evaluated the splicing status of these transcripts by  
8 analyzing three parameters: the genomic distance between the L1 side and the unique sequence for  
9 each chimera, the total chimera size, and the position of the splice junctions identified by TopHat  
10 around the L1 first nucleotide in the 11 samples (Figure S4). We found evidence of splicing for 44  
11 LCTs (including for the 17 spliced LCT-ESTs already described). This indicated that most LCTs  
12 identified by CLIFinder resulted from continuous transcription from the L1-ASP to the unique  
13 adjacent sequence and that only 2.9% of LCTs were spliced at their 5'end.

14 Our strategy, based on Ribo-Zero RNA-seq, identifies LCTs regardless of their polyadenylation  
15 (polyA) status. To determine whether LCTs were polyA transcripts, we tested the presence of 24  
16 LCTs, selected among the 31 we previously validated, in cDNA obtained by reverse transcription of  
17 total RNA from two glioma samples (GS#1 and #2) with random hexamers or with an oligodT  
18 oligonucleotide. We could amplify all 24 LCTs in cDNA obtained with random hexamers, but only 11  
19 in cDNA reverse transcribed with the oligodT primer (Fig. 2a). The absence of amplification for 13  
20 LCTs could be due to the presence of a non-polyA transcript or by inefficient full-length reverse  
21 transcription (*e.g.* long transcript). Therefore, we incubated the same RNA samples with a poly(U)  
22 polymerase to add an artificial polyU tail to their 3'end, before the reverse transcription step  
23 performed using an oligodA oligonucleotide. In this case, we could amplify 22 of the 24 LCTs (Fig.  
24 2a), suggesting that at least 11, and up to 13 of the 24 studied LCTs were not polyadenylated.

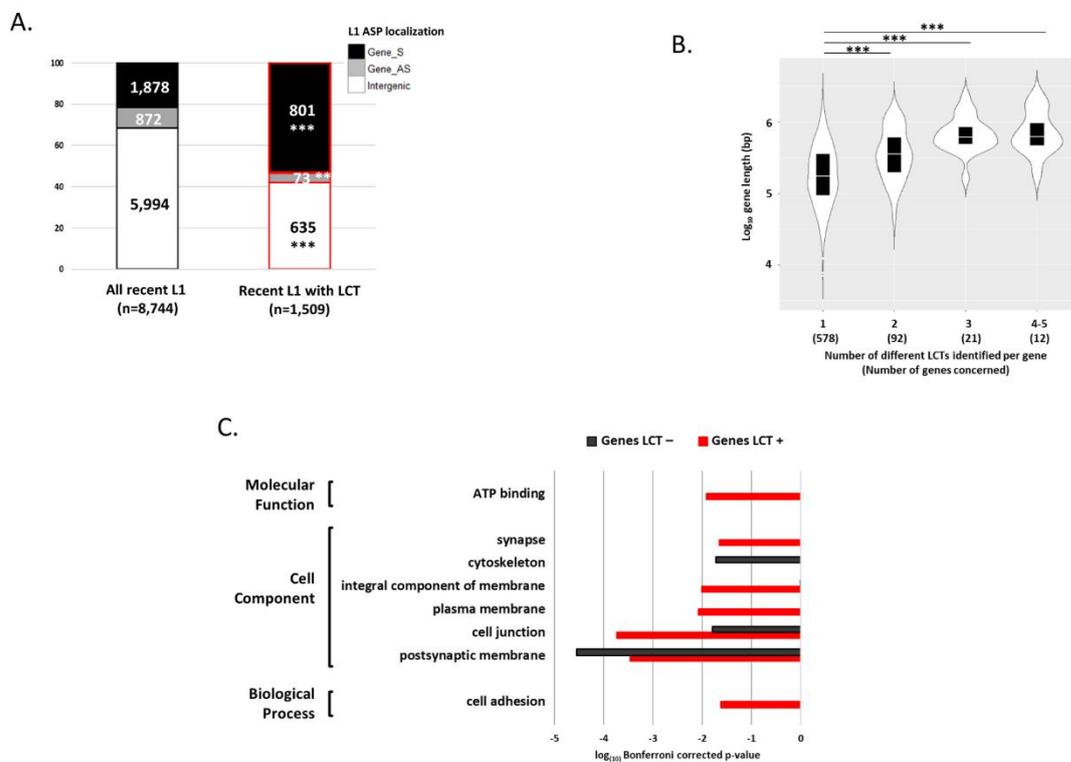
25 In a more systematic approach, we compared the data obtained by Ribo-Zero RNA-seq of GSC1 and  
26 GSC2 samples with those obtained by polyA RNA-seq of GSC2 and GSC6 samples. The number of  
27 LCTs detected per sample was always higher in the Ribo-Zero RNA-seq datasets (Fig. 2b). Similarly,



1 **LTC-producing L1 retrotransposons in brain and *IDW*wt glioma samples tend to localize in**  
 2 **genes with brain-related functions**

3 Compared with the 8,744 recent L1 retrotransposons possessing a 5' UTR, the 1,509 L1  
 4 retrotransposons involved in LCT production did not show a specific size (n=1,115, 73.9%, were 6kb  
 5 in length; data not shown) or genomic distribution (Figure S5).

6 Conversely, 58% of them (874/1,509) were intragenic, compared with 31.5% for all recent L1 in the  
 7 human genome (Fig. 3a). Surprisingly, only 8% of intragenic L1-ASPs producing LCTs were in the  
 8 antisense orientation relative to their host gene, compared with 32% for the whole L1 population  
 9 considered here. Most host genes (n=578) contained only 1 LCT-producing L1, 92 genes contained 2  
 10 LCT-producing L1, and 33 genes contained 3 to 5 LCT-producing L1 retrotransposons. The latter  
 11 genes were large genes with a size ranging from 0.16 to 2.3 Mb (Fig. 3b). Gene Ontology analysis  
 12 showed that the host genes of the 874 intragenic LCT-producing L1 were significantly enriched in  
 13 genes involved in brain function, *i.e.* synapse function and components (Fig. 3c). The host genes of  
 14 recent L1 elements that do not produce LCTs in brain samples were enriched in genes with more  
 15 ubiquitous functions, such as cytoskeleton components.

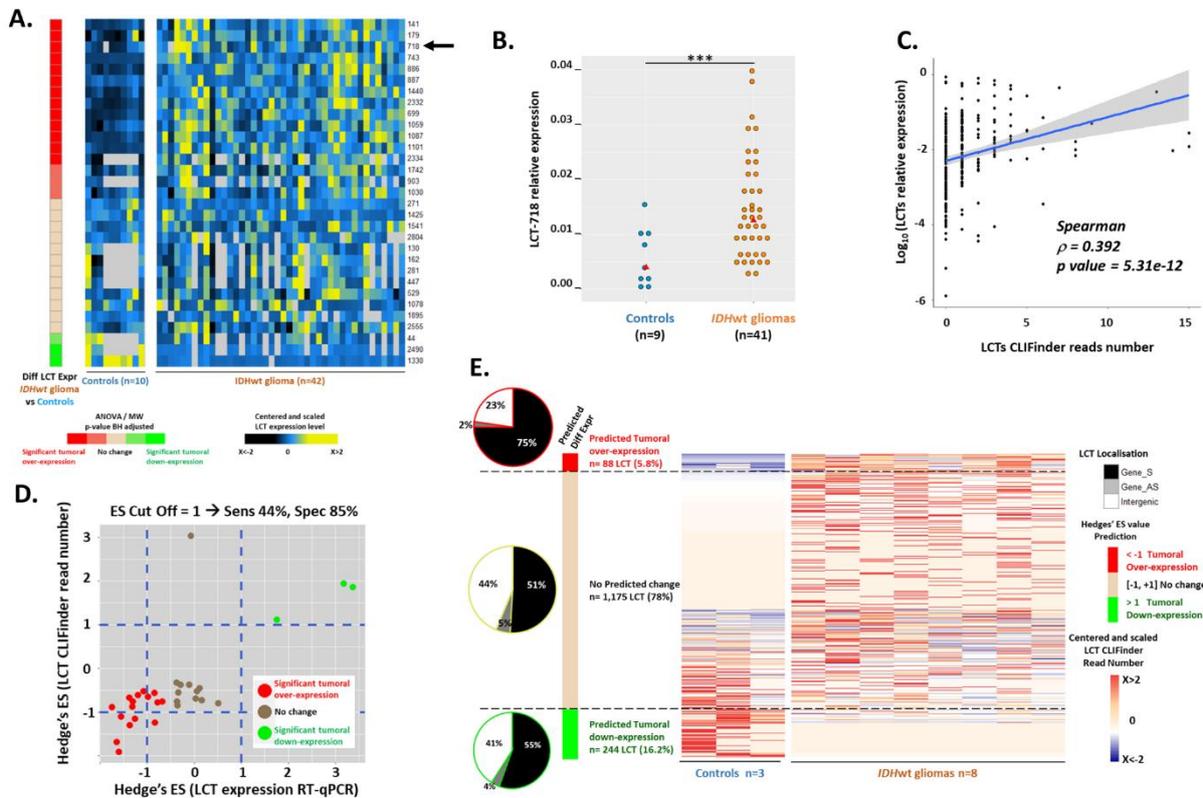


1 **Figure 3: Characteristics of the 1,509 LCTs identified by CLIFinder in *IDHwt* glioma and**  
 2 **normal brain samples**

3

4 **LCTs are expressed in healthy brain and aggressive glioma**

5 As previous studies suggested that some LCTs could be tumor-specific (7,19,20), we asked whether  
 6 LCT expression differed in glioma and healthy brain samples. As shown before (Figure S1), extensive  
 7 genome-wide LCT identification required to pool data of several samples. In agreement, CLIFinder  
 8 analysis of the RNA-seq data from three healthy brain and eight *IDHwt* glioma samples did not allow  
 9 assessing the LCT transcriptional status in these two tissue types. Therefore, we first quantified, by  
 10 RT-PCR, the expression of the 31 validated LCTs in 10 brain control and 41 *IDHwt* glioma samples.  
 11 All 31 LCTs were expressed in both control and glioma samples (Fig. 4a), although at different levels.  
 12 Specifically, 16 were upregulated (*e.g.* LCT 718; Fig. 4b) and 3 were downregulated in glioma  
 13 samples (Fig. 4a).



14

15 **Figure 4: LCTs are expressed in normal brain and glioma samples, but with different expression**  
 16 **levels for a LCT subset.**

1  
2 Moreover, the relative expression, measured by RT-qPCR, of the 31 validated LCTs was positively  
3 correlated with their reads number determined by CLIFinder (Spearman's  $Rho = 0.395$ ,  $p < 2.2e-16$ )  
4 (Fig. 4c). This suggests that the CLIFinder reads number represents a semi-quantitative value than can  
5 be used to predict the transcriptional status of the 1,509 LCTs in brain and tumor samples. Therefore,  
6 we calculated the Hedges'  $g$  effect size (Hedges' ES) (28). To establish a Hedge's  $g$  cut-off value that  
7 represented the best compromise between sensitivity and specificity, we compared the  $g$  values of the  
8 31 validated LCTs, obtained using their RNA-seq reads number, with those obtained using their  
9 relative expression by RT-qPCR and annotated according to their differential expression in tumor  
10 samples compared with controls (Fig. 4d). This allowed us to define a Hedges'  $g$  cut-off value of  $\pm 1$   
11 (44% of sensitivity and 85% of specificity) (Fig. 4d).

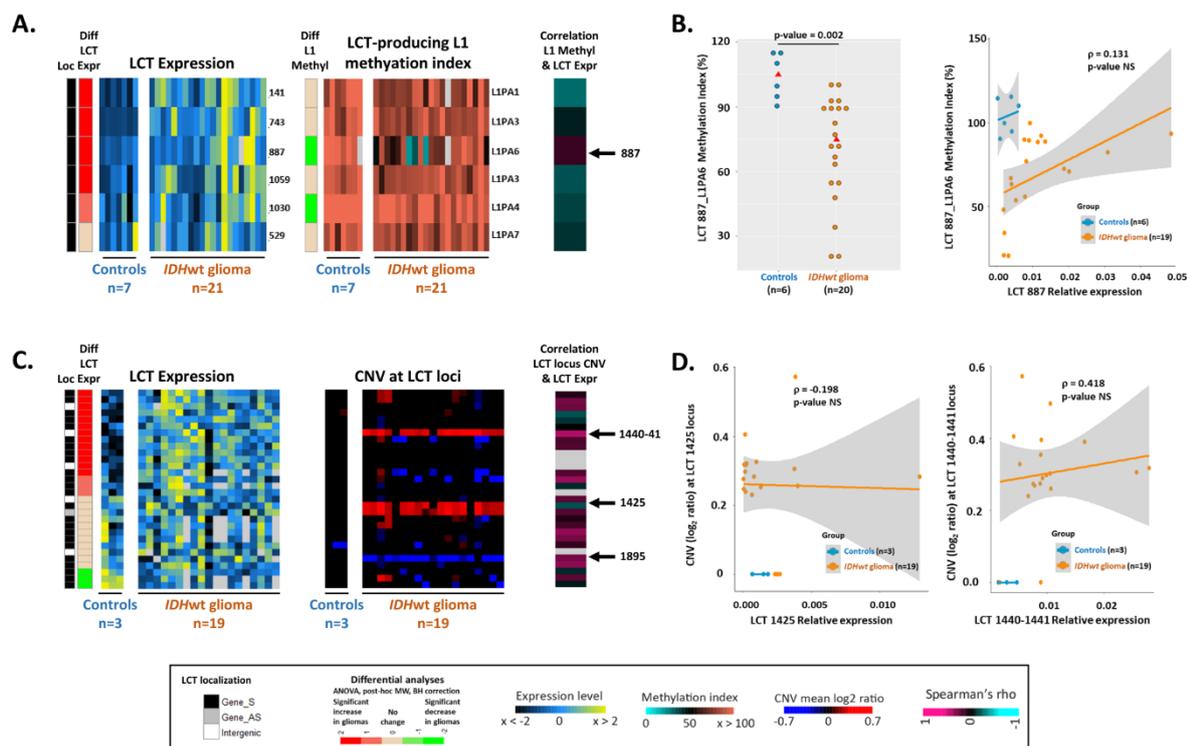
12 Using this cut-off value, the Hedges'  $g$  values obtained for the 1,509 LCTs predicted that 244 (16.2%)  
13 and 88 (5.8%) of them were down- and up-regulated, respectively, in *IDHwt* glioma *versus* control  
14 samples ( $g$  values:  $>1$  and  $<-1$ , respectively) (Fig. 4e). Note that the "upregulated" category does not  
15 distinguish between LCTs that are overexpressed and LCTs that are expressed only in glioma samples.  
16 Altogether, these observations indicated that most LCTs are not glioma-specific. They suggested that  
17 most LCTs are produced in healthy brain and also aggressive glioma, and that approximately 20% of  
18 them shows tumor-specific down- or up-regulation.

19

## 20 **Aberrant DNA methylation and genome copy number variations are not the main cause of LCT** 21 **expression deregulation in glioma samples**

22 To elucidate the causes of the tumor-specific upregulation of some LCTs, we first asked whether it  
23 could be explained by L1 promoter hypomethylation, resulting from the global hypomethylation of  
24 tumor DNA, as previously demonstrated for the *L1-MET* LCT in cancer cell lines (17). To test this  
25 hypothesis, we optimized the qAMP technique (29) to assess the promoter methylation level in 6 of  
26 the 31 validated LCT-producing L1 elements. Five of the associated LCTs were overexpressed in  
27 glioma samples, while one displayed the same expression level in control and glioma samples. The

1 qAMP approach, based on the use of methylation-sensitive and methylation-dependent restriction  
 2 enzymes followed by real-time PCR, allows the analysis of several CpG sites localized in the  
 3 immediate 5' L1 sequence, thus covering the CpG sites previously described as informative (17)  
 4 (Figure S6 a-c). Comparison of the methylation index obtained for each of the six loci in 7 brain  
 5 controls and 21 *IDHwt* glioma samples identified two L1 elements associated with LCTs  
 6 overexpressed in tumors and in which promoter methylation was reduced in the *IDHwt* samples  
 7 compared with controls (-28% for LCT-1030,  $p=0.001$ ; and -30% for LCT-887,  $p=0.002$ ) (Fig. 5a-b;  
 8 Figure S6c). However, we did not detect any significant negative correlation between LCT relative  
 9 expression and promoter methylation of the associated L1 element at these two hypomethylated loci  
 10 nor at the other four. Moreover, the methylation index for the six loci remained high in most glioma  
 11 samples.



12  
 13 **Figure 5: L1 hypomethylation and CNVs are not the driving force of LCT expression changes in**  
 14 ***IDHwt* glioma samples**

15

1 Next, to determine whether genomic rearrangements, which are often observed in cancer cells, could  
2 contribute to LCT deregulation, we analyzed Copy Number Variation (CNV) in the genomic loci of  
3 the 31 validated LCTs in 3 controls and 19 *IDHwt* samples. *IDHwt* samples are characterized by  
4 chromosome 7 gain and chromosome 10 loss (24). Accordingly, most *IDHwt* glioma samples carried  
5 extra and fewer copies of the L1 elements located on chromosome 7 and 10, respectively (Fig. 5c).  
6 However, there was no obvious link between CNV in glioma samples and LCT expression variation.  
7 Indeed, among the three LCT-producing L1 elements on chromosome 7 (LCT\_1425, LCT\_1440-41  
8 and LCT\_1451), only LCT\_1440-41 was overexpressed in glioma samples. Moreover, expression of  
9 the chromosome 10-associated LCT\_1895 was unchanged in glioma samples relative to controls,  
10 although the L1 copy number was decreased in glioma samples. Finally, the 15 loci associated with  
11 the other overexpressed LCTs did not show any CNV. Moreover, comparison of CNV and expression  
12 data at each LCT in each sample did not highlight any significant correlation between CNV and LCT  
13 expression (Fig. 5c and d).  
14 Combined, these observations indicated that aberrant DNA methylation and genome CNV might  
15 contribute, but are not the main cause of LCT deregulation in glioma samples.

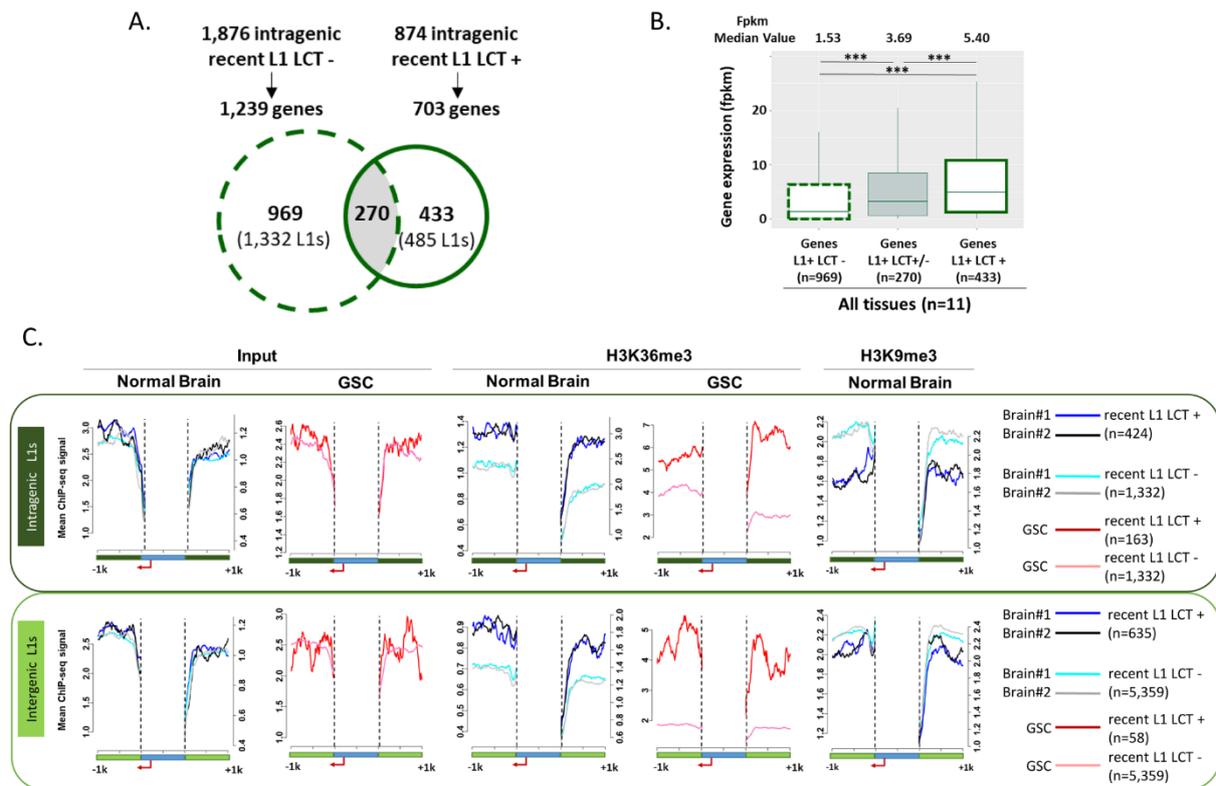
16

### 17 **LCT transcription from the L1-ASP occurs in transcriptionally active loci**

18 The previous data argue against a central role for DNA methylation in LCT production regulation.  
19 Specifically, the observation that LCT-producing L1 elements in control and *IDHwt* glioma samples  
20 tend to localize in genes with brain functions suggests that the L1-ASP activity of recent L1 could be  
21 influenced by the host locus transcriptional activity, as observed for the sense promoter (30,31). To  
22 test this hypothesis, we first compared the expression level of genes that contain one or several LCT-  
23 producing L1 retrotransposons (n=433 genes/485 L1) and of genes hosting one or several recent L1  
24 not associated with LCT production (n=969 genes/1,332 L1) (Fig. 6a). Comparison of the RNA-seq  
25 data from normal brain (n=3) and *IDHwt* glioma (n=8) samples indicated that the expression level of  
26 genes containing LCT-producing L1 retrotransposons was significantly higher (Fig. 6b). This was true  
27 when the analysis included all samples (control and glioma) (Fig. 6b), and also when focused only on

1 controls or aggressive glioma samples (Figure S7a). Thus, intragenic LCT transcription from the ASP  
 2 for recent L1 retrotransposons seems to occur preferentially from transcriptionally active genes.

3



4

5 **Figure 6: Recent LCT-producing L1 retrotransposons are localized in transcriptionally active**  
 6 **regions**

7

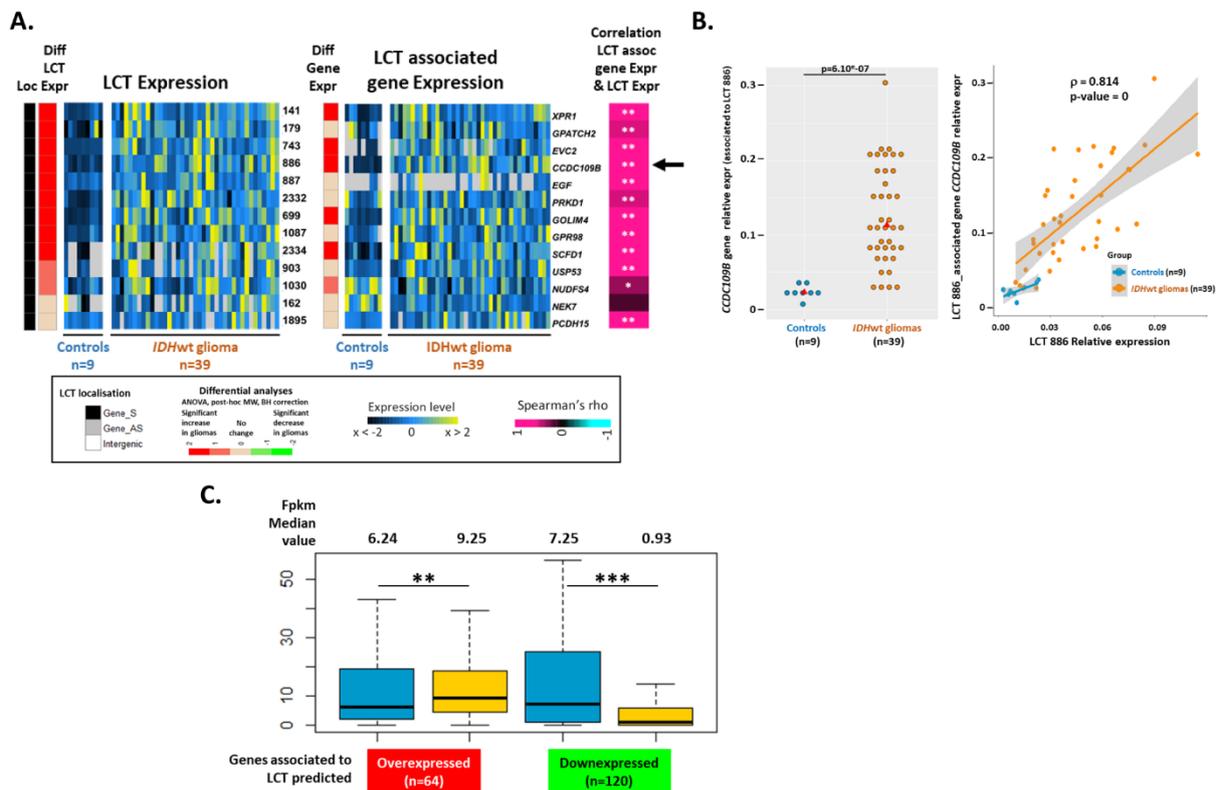
8 In agreement, analysis of chromatin data of the -1kb/+1kb regions surrounding recent L1  
 9 retrotransposons from two normal brain samples and one GSC line highlighted that the genomic  
 10 regions surrounding LCT-producing L1 elements were enriched in histone marks associated with  
 11 transcribed regions (*i.e.* H3K36me3) compared with those of recent L1 elements not associated with  
 12 LCT production (Fig. 6c, and Figure S7b). Conversely, the repressive H3K9me3 mark was enriched  
 13 around regions of recent L1 retrotransposons not associated with LCT transcription (Fig. 6c). We  
 14 obtained similar results for intergenic L1 retrotransposons, suggesting that those LCT-producing L1  
 15 elements are localized in transcribed intergenic regions (Fig. 6c).

1 Altogether, these observations support the hypothesis that LCT transcription from inter- and intra-  
 2 genic L1 retrotransposons in normal brain and *IDHwt* glioma samples is associated with the active  
 3 transcriptional state of the host locus. Like for the L1 sense promoter, the transcriptional activity of the  
 4 L1-ASP might be influenced by its immediate surrounding genomic sequence.

5

## 6 LCT expression deregulation in tumors follows that of their host locus

7 Then, to determine whether the tumor-associated change in the expression of some LCTs could be  
 8 explained by transcriptional deregulation of the host locus, we assessed by RT-qPCR the expression of  
 9 13 genes hosting 13 validated LCTs (11 overexpressed and 2 unchanged in tumors) in 9 controls and  
 10 39 *IDHwt* glioma samples (Fig. 7a). Six genes were significantly overexpressed in tumor samples  
 11 (Fig. 7a-b). Moreover, correlation analyses indicated that for all host loci with overexpressed LCTs,  
 12 the host gene expression was highest in samples with the highest LCT expression (Fig. 7a-b).  
 13 Univariate linear regression analysis validated these correlations, and multivariate linear regression  
 14 analysis confirmed the strong positive correlation between the expression of the LCT and associated  
 15 gene, independently of the sample group.



16

17

1 **Figure 7: LCT deregulation in glioma samples is linked to transcriptional changes in the host**  
2 **gene**

3

4 We next asked whether this observation could be applied to all intragenic LCTs. Therefore, we  
5 compared the expression (in fpkm, from RNA-seq data) of genes associated with the 64 and 120  
6 intragenic LCTs predicted to be up- and down-regulated in gliomas, respectively, in controls and  
7 *IDHwt* glioma samples. The gene expression in tumor and control samples was similar to that of their  
8 associated LCTs (Fig. 7c). Genes associated with upregulated LCTs also were significantly  
9 overexpressed in *IDHwt* glioma samples, whereas genes associated with downregulated LCTs also  
10 were significantly downregulated.

11 This suggests that LCT down- and up-regulation in glioma imply a similar transcriptional deregulation  
12 of the relevant host locus.

13

14 **Discussion**

15 Our study showed that LCT transcription from the ASP region in the L1 5'UTR is widespread in  
16 aggressive glioma, and that 17.25% (n=1,509) of recent L1 retrotransposons are involved in their  
17 production. It also revealed that these events are not tumor-specific, because LCT transcription  
18 occurred also in healthy brain, as previously shown in normal colon (32).

19 To date, *in silico* and molecular analyses in various human normal tissue and cancer samples and cell  
20 lines identified 161 recent LCT-producing L1 retrotransposons (4,7,14,17,18,20). Here, by analyzing  
21 healthy brain and *IDHwt* glioma samples using Ribo-Zero RNA-seq and a dedicated bioinformatics  
22 tool, CLIFinder (23), we extended the LCT landscape by revealing that more than 1,500 recent L1  
23 retrotransposons can produce LCTs in both sample types. We also identified two LCT features that  
24 mainly explain this marked increase in the ability to detect LCTs: most are unspliced at their 5' end,  
25 and they can be polyA or non-poly-A. Unlike previous approaches (4,7,14,17,18,20), our strategy is  
26 not biased toward spliced transcripts and uses total RNA. Therefore, this strategy requires RNA-seq  
27 data with high reading depth and/or the pooling of data from several samples to be representative. The

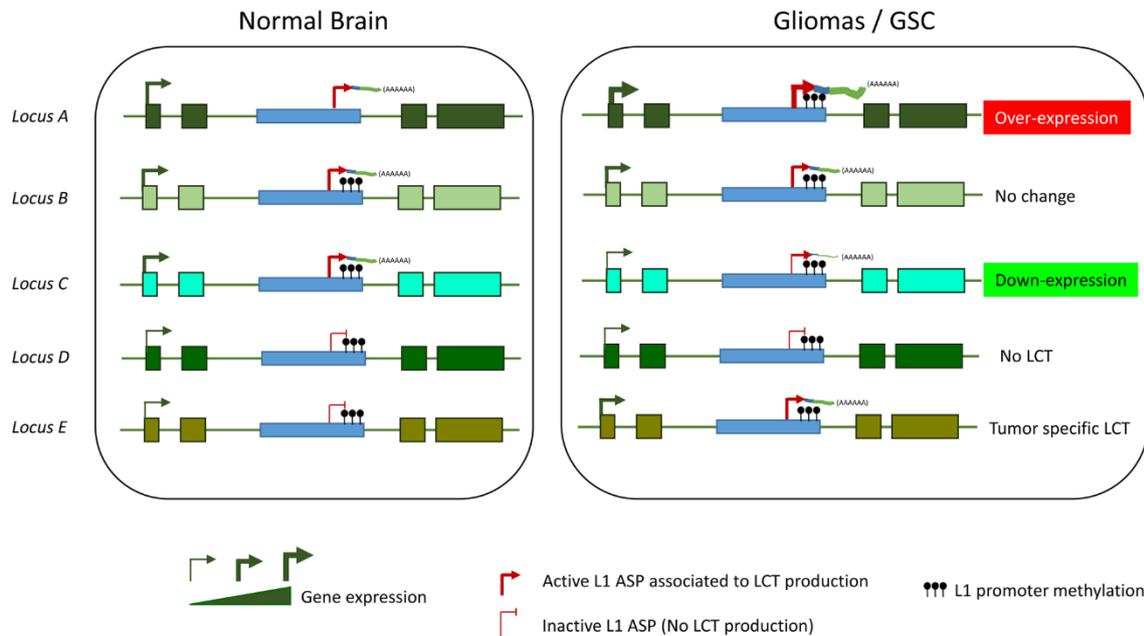
1 vast majority of LCTs we identified are unspliced at their 5' end, although two splicing donor sites are  
2 located in the L1 promoter. This finding was obtained using RNA-seq data based on total RNA and  
3 also polyA RNA, arguing against a technical bias linked to nascent RNA sequencing, and stresses that  
4 continuous transcription from the L1-ASP to the adjacent unique sequence is a major feature of the  
5 LCT 5' end. However, this does not preclude the possibility of downstream splicing events in the LCT  
6 transcripts, as observed in most of the previously characterized LCTs. Moreover, 89% of the LCTs we  
7 identified in GSC lines are also found in glioma samples isolated from different patients,  
8 demonstrating that our approach allows revealing a LCT tumor-type specific signature. Altogether,  
9 these observations validate our approach as relevant to investigate genome-wide LCT expression from  
10 L1-ASP in a tissue type. It should be noted that due to the CLIFinder settings used here, our  
11 conclusions are valid only for LCTs that initiate at ASP regions located in the first 500bp of the L1  
12 elements. However, our analysis of long-read sequencing datasets revealed that LCTs may also initiate  
13 in the coding sequence of L1. To better characterize these LCTs and to determine whether they may  
14 have arisen from older subfamilies, such as L1PA8, they should be compared with a reference set of  
15 L1 coding sequences (*i.e.* 900-1900 bp). The tunable nature of CLIFinder allows performing other  
16 refined analysis. For instance, it can be used to identify also a tumor sample-specific, rather than  
17 tumor type-specific signature. Such tumor type-specific analyses can be necessary for taking into  
18 account potential tumor-related genome rearrangements and the presence of additional L1 copies not  
19 present in the reference genome. In this case, the tumor genome, if available, should be used for the  
20 CLIFinder alignment step, and not the human reference genome.

21 A key finding of our study is that the L1-ASP transcriptional activity and that of the host locus are  
22 coupled (Fig. 8). This observation provides an experimental evidence to the model proposed by  
23 Nigumann *et al.* (2002) (14) about two decades ago according to which L1-ASP activation “*could be a*  
24 *consequence of the activation or chromatin opening of the cellular host genes*”. In agreement, we  
25 observed that in normal brain and in GSC lines, LCT-producing L1 retrotransposons are embedded in  
26 genomic regions enriched in H3K36me3, a chromatin mark associated with transcription elongation.  
27 As intergenic LCT-producing L1 retrotransposons present a similar feature, we propose that  
28 transcription events in general, regardless of their occurrence in intragenic or intergenic regions, can

1 influence the L1-ASP promoter activity. This is similar to what described for intragenic L1 sense  
2 promoter activation that is restricted to the cells in which the host gene is transcribed (30,31),  
3 indicating that the host locus transcriptional activity influences the whole L1 5'UTR transcriptional  
4 activity. However, as sense and antisense L1 promoter activities generally do not co-exist on the same  
5 L1 sequence (6,31), other regulation levels are likely to refine transcription from the sense and  
6 antisense promoters, particularly the relative orientation of the L1 element and the host genes. Indeed,  
7 68% of all intragenic recent L1 retrotransposons are in the opposite orientation as their host gene, and  
8 this percentage increases to 92% for LCT-producing L1 retrotransposons in brain and glioma samples.  
9 This suggests that sense-oriented progression of the transcriptional complex through the L1-ASP (or  
10 L1 sense promoter) is required to promote the transcriptional activity and/or the stability of the  
11 resulting LCT (or L1) transcript. This model, the underlying molecular bases of which remains to be  
12 determined, provides a mechanism whereby L1 sense and antisense promoter activities are mutually  
13 exclusive. It also ensures that intragenic L1 retrotransposons are not a source of antisense transcripts  
14 relative to those of the host gene that could affect the locus epigenetic regulation and mRNA stability  
15 through the formation of double-stranded RNA. Our study also bring some insights into the role of  
16 DNA methylation in the control of L1-ASP activity. In line with a comprehensive study on two L1-  
17 ASPs in various cancer cell lines (33), but unlike what described for the L1-*MET* LCT (17), we did not  
18 observe any inverse correlation between L1 promoter methylation and L1-ASP activity (evaluated by  
19 measuring the relative LCT expression level). It must be noted that in all studies on this question  
20 (17,33), including the present work, DNA methylation was assessed at CpG sites located in a CpG  
21 island that overlaps with the internal L1 sense promoter region, but downstream of the L1-ASP  
22 transcription starting site. Therefore, the role of this CpG island in the control of the L1-ASP could be  
23 questioned. Moreover, the L1-ASP core activity and putative regulatory regions (position +450 to  
24 +800 of the 5'UTR) (4) are depleted in CpG sites, further arguing against a central role of DNA  
25 methylation in the control of L1-ASP activity. Altogether, these data suggest that the host gene locus  
26 transcriptional activity is the main driving force of LCT expression in physiological and pathological  
27 conditions. Therefore, and unlike what observed at the L1 sense promoter, L1-ASP activity is more  
28 influenced by the cancer-associated changes in the host gene expression level rather than by the  
29 cancer-associated disruption of the DNA methylation pattern. This leads to the question of whether

1 and how LCTs play a functional role during tumorigenesis. In a simple model, tumor-specific LCTs,  
 2 hosted by genes expressed in the tumor, might influence the expression of their host and/or  
 3 surrounding genes, thus directly contributing to tumorigenesis. However, in aggressive glioma, this  
 4 model is not valid for the majority of LCTs because they are also expressed in healthy brain.

5



6

7 **Figure 8: Working model of LCT transcriptional regulation in normal brain and in**  
 8 **glioma/GSCs**

9

10 We could hypothesize that changes in LCT expression at a given locus might modify its role in cancer  
 11 cells compared with normal cells. *LI-MET* and LCT13, the two main examples of LCT role in  
 12 tumorigenesis, are both expressed in healthy cells. Their ability to promote oncogenic functions, as  
 13 proposed for *LI-MET* (17), or to affect the expression pattern of surrounding genes, as proposed for  
 14 LCT13 (19), is thought to be related to a change of their expression level in cancer cells (17,19).  
 15 Another open question is whether LCT expression level changes can modify their length and splicing

1 signature in tumor cells. Future studies should focus on this key question and on evaluating the coding  
2 potential of these LCT transcripts. Within this frame, we have initiated a long read sequencing  
3 approach to analyze three GSC samples (see figure 1e). This yet in development approach requires now  
4 to be conducted to more samples and to benefit from the development and validation of dedicated  
5 bioinformatics tools to identify full-length LCT sequences in healthy brain and aggressive glioma.

6 In conclusion, our study describes the full LCT expression profile in brain and glioma samples.  
7 Unexpectedly, this analysis highlighted that LCT expression is not restricted to glioma, but is  
8 widespread also in normal brain: approximately 17% of all recent L1 retrotransposons with a 5'UTR,  
9 from the L1PA1 to L1PA7 subfamilies, produces a LCT in brain and glioma. Our study also revealed  
10 that the transcriptional activity from L1-ASP and from the promoter of the host locus are coupled,  
11 suggesting that LCT transcription is mainly regulated by the host locus transcriptional activity and  
12 chromatin signature.

13

## 1 **Materials and Methods**

### 2 **Tumor and control samples**

3 Diffuse glioma samples from adult patients who underwent surgical resection between 2007 and 2014  
4 were obtained from Clermont-Ferrand University Hospital Center, France (“Tumorothèque Auvergne  
5 Gliomes”, ethical approval DC-2012-1584). This study was approved by the relevant ethics  
6 committees and competent authorities, and the study protocols follow the World Medical Association  
7 Declaration of Helsinki. Written informed consents were provided by all patients. Samples were  
8 isolated as previously described (34,35). In this study, the 42 included aggressive glioma samples  
9 carried wild type *isocitrate dehydrogenase (IDH)* gene (*IDHwt*) (36).

10 Ten control brain samples (healthy controls; samples removed by autopsy 4-16h after accidental death)  
11 were obtained from the Brain and Tissue Bank of Maryland (mean age of 27.3 years, standard  
12 deviation  $\pm$  2 years). These samples, identified by the Brain and Tissue Bank of Maryland as corpus  
13 callosum (n=8) and frontal cortex (n=7), correspond to white matter enriched in astrocytes and  
14 oligodendrocytes, and are relevant non-cancer controls for gliomas. Before use, tumor and control  
15 samples were homogenized by cryogenic grinding, and each sample was aliquoted in at least three  
16 vials for genomic DNA, RNA, and chromatin extraction. All samples were stored at -80°C until use.

17 Cell pellets from three GSC lines (GSC-1, GSC-2, and GSC-6) derived from patients with *IDHwt*  
18 gliomas were obtained from Poitiers University Hospital Centre, France, and were previously  
19 characterized (37–39).

### 20 **List of ASP-containing L1 elements**

21 The coordinates of L1 elements were obtained from the RepeatMasker database deposited in the  
22 UCSC Table Browser. Recent L1 elements were selected based on their RepName (L1P1, L1HS,  
23 L1PA2, L1PA3, L1PA4, L1PA5, L1PA6, L1PA7, L1PA8). This list was then filtered using repLeft  
24 and repStart <400 and repEnd >600 to retain only L1 elements with an ASP. The final list is in Table  
25 S1.

26

1

## 2 **Illumina RNA-sequencing and LCT identification using CLIFinder**

3 Total RNA was isolated from frozen tissue samples and frozen cell pellets as previously described  
4 (35). Strand-oriented RNA-seq was performed using total RNA (n=3 brain control, n=8 *IDHwt*, and  
5 n=2 GSC samples: GSC-1 and GSC-2) and also polyA mRNA from the GSC-6 and GSC-2 samples  
6 (34).

7 CLIFinder was used as previously described (23) to extract chimeras from stranded paired-end RNA-  
8 seq data. The analyses were based on the reference annotated human genome GRCh37/hg19. The  
9 reference set of L1 sequences used was a text file containing FASTA sequences that corresponded to  
10 the first 500 bp of the 9,123 L1 elements (L1PA1 to L1PA8) with a 5'UTR sequence in the human  
11 genome. These sequences were extracted from the Repeat Masker database with UCSC tools. As the  
12 reference file contains all potential targeted sequences, the parameters used to select chimeric cDNA  
13 of potential interest were:

- 14 - 5' end read sequence in which at least 50 bp that matched a L1 reference sequence (1  
15 mismatch was tolerated to take into account single nucleotide polymorphisms that may be  
16 present also in L1 sequences)
- 17 - the associated 3'end read sequence should contain at least one unique sequence (*i.e.* not  
18 annotated by Repeat Masker) of  $\geq 30$  bp.

19 Finally, a maximum insert size of 50 kb after alignment to genomic DNA was tolerated between the  
20 paired-end reads retained (to allow the identification of spliced chimeras).

21 In the 11 RNA-seq datasets, CLIFinder identified 1,677 chimeras associated with the L1 subfamilies  
22 L1PA1 to L1PA8 (Table S2). Among these chimeras, 37 were transcribed in the same orientation as  
23 the associated L1, and 69 were associated with L1 elements without conserved 5'UTR, therefore  
24 resulting in 1,571 chimeras that potentially corresponded to LCTs.

## 25 **Long-read SMRT sequencing and LCT TSS identification in GSC samples**

1 To sequence full-length polyA and non-polyA LCT transcripts, a modified Iso-seq protocol (PacBio)  
2 was used. For each sample, 3 µg of total RNA was poly-uridynylated using a polyU polymerase (New  
3 England Biolabs) following the manufacturer's recommendations. After poly-uridynylated (polyU)  
4 RNA purification with the RNeasy kit (Qiagen) and ethanol precipitation, two reverse transcription  
5 reactions were performed (each with 1µg of polyU RNA) using the SMARTer Kit (Takara) in which  
6 the oligodT<sub>(30)</sub> 3' SMART CDS Primer II A was replaced by the oligodA<sub>(30)</sub> 3' SMART CDS Primer II  
7 A. The obtained full-length cDNA (1:10 dilution) was then amplified by PCR (14 cycles) using GXL  
8 PrimeSTAR DNA polymerase (Takara) in 20 PCR reactions. PCR products were purified and  
9 concentrated (in 21µL) with two rounds of 1X AMPurePB beads (PacBio) treatment. cDNA size  
10 fractionation was then performed on 0.75% agarose gels (without ethidium bromide) to individualize  
11 three cDNA fractions: [0.3 - 2 kb], [2 - 5 kb], and [> 5 kb]. Dedicated biotinylated riboprobes,  
12 designed against the 5' 400 bp sequences of the 8,744 recent L1 (myBaits, Arbor Biosciences), were  
13 used to perform three independent captures of each cDNA fraction according to the supplier's  
14 recommendations (except for the addition of Block C solution = human Cot-1 DNA). Captured  
15 cDNAs were then amplified by PCR using Prime StarGXL DNA polymerase (Takara) with 10, 12 and  
16 14 cycles respectively for the <2 kb, 2-5 kb and >5 kb cDNA fractions. PCR products were pooled  
17 and purified using 0.55X AMPurePB beads (PacBio). The three samples were pooled together and  
18 handled for SMRTbell template preparation and sequencing (Sequel II, Gentyane platform). The  
19 obtained sequences were treated as follows: 1) consensus sequence establishment using the Iso-seq3  
20 CCS algorithm; 2) selection of the full-length cDNA sequences including the SMARTer IIA adapter at  
21 both their 3' and 5' ends using the Iso-seq3 LIMA tool; 3) trimming of the polyT tail at the 3'end of  
22 each sequence (corresponding to the polyU tail added initially); and finally 4) alignment of the  
23 obtained sequences to the hg19 reference human genome using the minimap2 aligner. Only good  
24 quality (MapQ ≥30) aligned reads were retained in a SAM file (specifying chromosome, start and end  
25 genomic coordinates, transcription strand, MapQ value). The genomic coordinates of PacBio read  
26 alignments were extracted using the bamtoBED tools from the BEDtools suite V2.30.0 and the SAM file.  
27 Using the GenomicRanges R library, these coordinates were filtered to keep only reads that  
28 overlapped with the 226 L1-producing LCTs identified by CLIFinder in GSC Illumina RNA-seq  
29 (GSC1 & GSC2) data. Only PacBio reads with potential LCT characteristics were selected: reads

1 containing a 5' sequence antisense to one of the 226 L1 and a unique genome sequence at its 3' end. To  
2 identify the TSS of the PacBio reads within the L1s, the starting position of each read was repositioned  
3 using the L1 coordinates as reference. These data were used to determine the number of L1 that were  
4 associated with a PacBio read that initiated in the ASP region or the 5'UTR of the L1 (i.e. position 0-  
5 500 and 0-900 of a L1). To visualize the TSS position of all reads inside each of the 226 L1 elements,  
6 the 5' initiation position of all reads was considered to compute the percentage of reads that initiated  
7 along each L1.

#### 8 **Whole gene expression (fpkm) determination**

9 Illumina stranded RNA-seq reads from three controls and eight *IDHwt* glioma samples were used to  
10 determine the normalized gene expression (fpkm), as previously described (34). Briefly, reads were  
11 mapped to the human genome (hg19) using TopHat2 (version 2.1.0) and a transcript annotation file  
12 from GENCODE (Release 19). Gene expression levels were determined with Cuffquant and Cuffnorm  
13 from the Cufflinks suite (version 2.2.1).

#### 14 **Copy Number Variation analyses**

15 CNV analyses were performed using the Genome-Wide Human CytoScan HD Array (Affymetrix) and  
16 3 controls and 19 *IDHwt* glioma samples, as previously described (34). To determine whether CNV  
17 could be correlated with LCT expression level, the mean log<sub>2</sub> ratio of the CNV values and the relative  
18 expression level for each LCT locus in the same sample were compared using the Spearman's  
19 correlation.

#### 20 **Data access**

21 Data are available at the NCBI Gene Expression Omnibus (GEO; <https://www.ncbi.nlm.nih.gov/geo/>)  
22 under the following accession numbers: GSE123892 for the stranded total RNA-seq of control  
23 samples and *IDHwt* gliomas; GSE161438 for the stranded total RNA-seq of GSCs; GSM4907328 and  
24 GSM180209 for the stranded polyA RNA-seq of GSCs, GSE190930 for the PacBio sequencing of  
25 GSCs, and GSE123682 and GSE161275 for the Cytoscan HD data.

#### 26 **Chimera validation**

1 The LCT-Unique Sequence (LCT-US) primers (using the Primer 3 software: <http://primer3.ut.ee/>)  
2 were designed for the 36 chimeras selected according to their detection by CLIFinder using high or  
3 low reads number, they tumor-specificity or not, and their localization (intragenic or intergenic). This  
4 primer was combined with a L1 common primer (localized at the beginning of the 5' sequence of L1  
5 elements). Due to sequence homology, a common primer could be designed for the L1PA1 to L1PA7  
6 subfamilies and a specific L1 common primer for the L1PA8 subfamily. Then, chimera detection was  
7 validated by RT-PCR using RNA from two glioma samples (GS#1 and #2) already used for RNA-seq.  
8 Briefly, 1 µg of total RNA was treated with DNase I (Promega) according to the manufacturer's  
9 recommendations, and then reverse transcribed using Random Hexamers (Invitrogen) and the  
10 SuperScript III Reverse Transcriptase Kit (Invitrogen). After dilution (1:5), 2 µl of cDNA solution was  
11 used for PCR amplification with the relevant LCT specific primer, LCT-US, and L1 common primers.  
12 Primer sequences, PCR conditions, and PCR product amplification size are given in Table S3. RT-  
13 PCR products were analyzed on 2% agarose gel and those of the expected size were cloned using the  
14 pGEM-T Easy Cloning Kit (Promega). One recombinant clone per chimera was sequenced to  
15 determine whether the obtained cDNA corresponded to the expected sequence.

#### 16 **RT-PCR-based 5' L1 walking approach to localize the TSS of the validated chimeras**

17 RNA available in large amounts and of good quality (RIN >8) from two glioma samples was used for  
18 these experiments. Additional L1 primers, downstream and upstream of the L1-ASP +450 location,  
19 were designed using the Primer3 software. Due to the high sequence homology displayed by the  
20 5'UTR sequence of the L1PA1 to L1PA6 subfamilies, the same L1 downstream and upstream primers  
21 were used for chimeras associated with L1 elements from these six subfamilies. However, as some  
22 mismatches could impair PCR amplification (when too numerous or when occurring at the 3'end of  
23 the primer), three different L1 downstream primers (PA1-PA6) were designed and alternatively used.  
24 Specific L1 downstream and upstream primers were designed for the L1PA7 and L1PA8 subfamilies.  
25 Primer sequences, mismatch number between L1 primer/L1 sequence, and PCR product amplification  
26 sizes and results are given in Table S3. As previously described, 1 µg of total RNA was treated with  
27 DNase I and reverse transcribed with SuperScript III using random hexamers. For each sample, a  
28 control was prepared without reverse transcriptase. PCR amplifications were performed as previously

1 described, using the relevant LCT-US primer and L1 downstream or upstream primer. In the PCR  
2 program, the extension time at 72°C was adjusted in function of the expected amplification product  
3 size (40 sec and 1 min for L1 downstream and L1 upstream primers, respectively). As all tested  
4 chimeras corresponded to continuous transcription events from the L1 sequence and the adjacent  
5 unique sequence, genomic DNA was used as positive control and as size marker. PCR products were  
6 analyzed on 1.5% agarose gel.

### 7 **Splicing analysis**

8 Three parameters were taken into account to identify splicing events in chimeras: 1) the distance  
9 between the start position of the unique sequence and the end position of the L1 sequence of each  
10 chimera after alignment (R1-R2 distance). A negative or null distance indicates an overlap or a  
11 juxtaposition of both sequences, and implies that these chimeras have been identified in multiple reads  
12 (either in one or different samples). A distance between 1 and 200 bp is concordant with the size  
13 selection made during the library construction (i.e. 280-340 bp) and often relevant of chimeras  
14 identified in few reads in all samples. A distance >500 bp may suggest at least one splicing event  
15 between the L1 and unique sequence; 2) the LCT size. Chimeras with a size >1000 bp, even if their  
16 R1-R2 distance is negative, can present evidences of splicing event(s), because the unique sequence is  
17 very large (> 700bp) due to concatenation of multiple reads by CLIFinder; 3) the occurrence of splice  
18 junctions identified by TopHat in the RNA-seq data in the  $\pm 200$  bp region around the L1 start  
19 position. One of these parameters was considered sufficient to conclude that a LCT was spliced. LCT  
20 with R1-R2 distance between 200 and 499 bp and a size between 600 and 999 bp were considered to  
21 have an undetermined spliced status.

### 22 **RT-PCR analysis to determine the polyA and non-polyA status of LCTs**

23 Total RNA from the glioma samples GS#1 and GS#2 was treated with DNase I and reverse  
24 transcribed with SuperScript III using 2.5  $\mu$ M oligod(T)<sub>n</sub> primer (Invitrogen) to specifically target  
25 polyA mRNA transcripts. Then, 2  $\mu$ L of diluted (1:10) cDNA was used for PCR amplification using  
26 the LCT-US and L1 common primers.

1 Alternatively, 3  $\mu\text{g}$  of total RNA was treated with DNase I and then with 2U of PolyU polymerase in  
2 the presence of 0.5mM UTP, 1X Buffer and 40U RNase inhibitor. The resulting poly-uridylated  
3 RNA was purified on RNeasy columns (Qiagen), concentrated after ethanol precipitation, and 1  $\mu\text{g}$  of  
4 PolyU-RNA was reverse transcribed using 0.6  $\mu\text{M}$  of oligod(A)n primer (Eurogentec) and Superscript  
5 III (Invitrogen). Then, 2  $\mu\text{L}$  of diluted (1:10) cDNAs was PCR amplified with the LCT-US and L1  
6 common primers.

### 7 **Expression analysis of LCTs and associated genes in gliomas by microfluidic RT-qPCR**

8 500 ng of total RNA was treated with DNase I (Promega) and divided in two independent aliquots  
9 (250 ng RNA/each) for reverse transcription using a random hexamer primer (Invitrogen). Then, first  
10 strand cDNA was pre-amplified (14 cycles) with the pool of primers used for the microfluidic PCR  
11 assays. Primer sequences and qPCR Efficiency (E) are given in Table S3. All primers used to quantify  
12 host gene expression are positioned upstream of the LCT TSS in L1. qPCR assays were performed and  
13 validated using Fluidigm 96.96 Dynamic Arrays and the Biomark HD system (Fluidigm) according to  
14 the manufacturer's instructions. LCTs and host genes were quantified in separate arrays. The relative  
15 expression level was quantified as  $R = (E_{\text{TOI}}^{-C_{\text{TOI}}}) / (\text{geometric mean } E_{\text{HK}}^{-C_{\text{HK}}})$  with E = PCR  
16 efficiency of each primer pair, TOI = Transcript Of Interest (*i.e.* chimera or gene to be quantified), and  
17 HK = housekeeping genes (*i.e.* *TBP*, *RPL13A* and *PPIA*) used to normalize transcript expression. For  
18 each sample, experiments were done in duplicate using the two independent RT reactions.

19 Differences in the expression levels of the LCT or associated gene between tumor and control samples  
20 were assessed with the one-way ANOVA and post-hoc Mann-Whitney test, if applicable. Significant  
21 p-values were adjusted with the Bonferroni correction. Correlation analyses between intragenic LCT  
22 and the associated gene expression levels were performed using Spearman's correlation.

### 23 **Chromatin analysis at recent L1 loci**

24 L1 positions and classifications were retrieved from the CLIFinder output. Based on these coordinates,  
25 chromatin analyses were done with the computeMatrix scale-regions tools from the deeptools suite  
26 (Version 3.1.3) (Parameters:--beforeRegionStartLength 1000, --regionBodyLength 1000, --  
27 afterRegionStartLength 1000, --skipZeros) and chromatin data coverages extracted from GeoDataset

1 in WIG format (for Brain input: GSM669971 and GSM773019; for GSC input: GSM1121871 and  
2 GSM1121861; for Brain H3K36me3: GSM670002 and GSM916041; for GSC H3K36me3:  
3 GSM1121868 and GSM1121858; for Brain H3K9me3: GSM916034 and GSM773017; for Brain  
4 H3K27ac: GSM773020 and GSM916035; for Brain H3K4me1: GSM669962 and GSM916039; for  
5 Brain H3K4me3: GSM670022 and GSM916040; for Brain H3K27me3: GSM916038 and  
6 GSM669913). These WIG files were converted into binary format with the wigToBigWig program  
7 from the UCSC server, to be compatible with the computeMatrix program prerequisites. Chromatin  
8 profiles surrounding the L1 start and end positions ( $\pm 1$  kb) that corresponded to the mean chromatin  
9 coverage over the set of genomic regions were generated using the computeMatrix output and R  
10 program.

### 11 **Quantitative analysis of DNA methylation by qPCR (qAMP)**

12 The qAMP approach relies on the use of methylation-sensitive (e.g. *HhaI* or *HpaII*) and methylation-  
13 dependent (e.g. *McrBC*) restriction enzymes. DNA was extracted from frozen samples (21 *IDHwt*  
14 glioma and 7 control brain samples) using the QIAamp DNA Mini Kit (Qiagen) according to the  
15 manufacturer's recommendations. For each DNA sample, four 200ng DNA aliquots were prepared.  
16 Three were incubated with 5U of *HhaI*, *HpaII*, or *McrBC* (NEB) at 37°C for 2h. The fourth aliquot  
17 was incubated only with buffer (negative control). The published qAMP protocol was slightly  
18 modified by adding a restriction enzyme inactivation step by incubation with proteinase K (0.4  
19 mg/ $\mu$ L) at 40°C for 30 min, followed by its inactivation at 95°C for 10 min. The obtained PCR  
20 templates were 8-fold diluted in water and stored at -20°C. PCR primers covered approximately the  
21 first 300 bp of each L1 promoter in which the restriction sites for each of the three enzymes used are  
22 present. L1 primers were designed using the Primer 3 software and combined to the already described  
23 LCT US primer. Additionally, primers for control regions with different DNA methylation levels were  
24 designed: 1) a control region that does not contain restriction sites for the three enzymes; 2) the  
25 promoter of the imprinted gene *PEG10* that is hemi-methylated; and 3) the *GAPDH* promoter region  
26 that in our samples, has a methylation level <10%. All primers used for the qAMP experiments are  
27 given in Table S3. Finally, artificially unmethylated (0% methylation) and fully methylated (100%)  
28 genomes were prepared by whole genome PCR amplification with the Repli-G Mini Kit (Qiagen) and

1 artificial methylation with the *SssI* enzyme (NEB). These two DNA templates were mixed to obtain a  
2 standard curve with theoretical DNA methylation levels of 0, 30, 70, and 100%. The efficiency of each  
3 primer pair was evaluated using 3  $\log_{10}$  gDNA concentrations. Only primer pairs with an efficiency  
4  $\geq 1.85$  and a standard curve corresponding to the expected theoretical DNA methylation percentages  
5 were retained. The DNA methylation index of each sample was then evaluated by qAMP on a  
6 LightCyclerR©480II (Roche). Four  $\mu\text{L}$  of each DNA (sham, *HpaII*-, *HhaI*- and *McrBC*-digested) was  
7 mixed with 1x SYBR Green Master Mix (Roche) and 0.5  $\mu\text{M}$  forward and reverse primers in a final  
8 volume of 10  $\mu\text{L}$ . The PCR program was: 95°C for 5 min, and then 40 cycles (95°C for 15 sec, 60°C  
9 for 15 sec, 72°C for 15 sec), and 72°C for 5 min. The DNA methylation index of each sample was  
10 calculated using the amplification data obtained for the three DNA samples digested with the  
11 restriction enzymes and as previously described using the  $\Delta\text{Ct}$  values ( $\text{Ct enzyme} - \text{Ct sham}$ ) (29).  
12 Differences in methylation index levels between glioma and control samples were assessed with the  
13 one-way ANOVA and post-hoc Mann-Whitney test, when applicable. Significant p-values were  
14 adjusted using the Bonferroni correction. Correlation analyses between L1 methylation index and LCT  
15 expression level were done using the Spearman's correlation.

16

1 **Acknowledgements**

2 We thank the Clermont-Ferrand hospital Neurosurgery department (Prof. J.J. Lemaire) for help in  
3 completing this study, the Platform “Gentyane” (<http://gentyane.clermont.inra.fr/>) for technical help  
4 with the Fluidigm 96.96 Dynamic Arrays and the SMRT Sequel II sequencing and the Bioinformatic  
5 Platform from GReD institute.

6 This study has been financed by grants from the Fondation ARC pour la Recherche Contre le Cancer  
7 (n° SFI20121205549) (to CVB), the Ligue Contre le Cancer Comité du Puy de Dôme (to CVB), the  
8 Plan-Cancer INSERM (CS14085)(to PA), the Fonds de Dotation Patrick Brou de Laurière (to PA and  
9 to CVB) and has been supported by the French government IDEX-ISITE initiative 16-IDEX 0001  
10 (CAP 20-25).

11 **Conflict of Interest statement**

12 The authors declare no competing or financial interests.

13

14

## 1 References

- 2 1. Lander, E.S., Linton, L.M., Birren, B., Nusbaum, C., Zody, M.C., Baldwin, J., Devon, K., Dewar,  
3 K., Doyle, M., FitzHugh, W., *et al.* (2001) Initial sequencing and analysis of the human genome.  
4 *Nature*, **409**, 860–921.
- 5 2. Brouha, B., Schustak, J., Badge, R.M., Lutz-Prigge, S., Farley, A.H., Moran, J.V. and Kazazian,  
6 H.H. (2003) Hot L1s account for the bulk of retrotransposition in the human population. *Proc.*  
7 *Natl. Acad. Sci.*, **100**, 5280–5285.
- 8 3. Maxwell, P.H. (2014) Consequences of ongoing retrotransposition in mammalian genomes. *Adv.*  
9 *Genomics Genet.*, **4**, 129–142.
- 10 4. Speek, M. (2001) Antisense Promoter of Human L1 Retrotransposon Drives Transcription of  
11 Adjacent Cellular Genes. *Mol. Cell. Biol.*, **21**, 1973–1985.
- 12 5. Khan, H. (2006) Molecular evolution and tempo of amplification of human LINE-1  
13 retrotransposons since the origin of primates. *Genome Res.*, **16**, 78–87.
- 14 6. Macia, A., Munoz-Lopez, M., Cortes, J.L., Hastings, R.K., Morell, S., Lucena-Aguilar, G.,  
15 Marchal, J.A., Badge, R.M. and Garcia-Perez, J.L. (2011) Epigenetic Control of Retrotransposon  
16 Expression in Human Embryonic Stem Cells. *Mol. Cell. Biol.*, **31**, 300–316.
- 17 7. Criscione, S.W., Theodosakis, N., Micevic, G., Cornish, T.C., Burns, K.H., Neretti, N. and Rodić,  
18 N. (2016) Genome-wide characterization of human L1 antisense promoter-driven transcripts.  
19 *BMC Genomics*, **17**, 463.
- 20 8. Denli, A.M., Narvaiza, I., Kerman, B.E., Pena, M., Benner, C., Marchetto, M.C.N., Diedrich, J.K.,  
21 Aslanian, A., Ma, J., Moresco, J.J., *et al.* (2015) Primate-specific ORF0 contributes to  
22 retrotransposon-mediated diversity. *Cell*, **163**, 583–593.
- 23 9. Lee, E., Iskow, R., Yang, L., Gokcumen, O., Haseley, P., Luquette, L.J., Lohr, J.G., Harris, C.C.,  
24 Ding, L., Wilson, R.K., *et al.* (2012) Landscape of Somatic Retrotransposition in Human Cancers.  
25 *Science*, **337**, 967–971.
- 26 10. Tubio, J.M.C., Li, Y., Ju, Y.S., Martincorena, I., Cooke, S.L., Tojo, M., Gundem, G., Pipinikas,  
27 C.P., Zamora, J., Raine, K., *et al.* (2014) Mobile DNA in cancer. Extensive transduction of  
28 nonrepetitive DNA mediated by L1 retrotransposition in cancer genomes. *Science*, **345**, 1251343.
- 29 11. Scott, E.C., Gardner, E.J., Masood, A., Chuang, N.T., Vertino, P.M. and Devine, S.E. (2016) A  
30 hot L1 retrotransposon evades somatic repression and initiates human colorectal cancer. *Genome*  
31 *Res.*, **26**, 745–755.
- 32 12. Scott, E.C. and Devine, S.E. (2017) The Role of Somatic L1 Retrotransposition in Human  
33 Cancers. *Viruses*, **9**.
- 34 13. Babaian, A. and Mager, D.L. (2016) Endogenous retroviral promoter exaptation in human cancer.  
35 *Mob. DNA*, **7**, 24.
- 36 14. Nigumann, P., Redik, K., Mätlik, K. and Speek, M. (2002) Many Human Genes Are Transcribed  
37 from the Antisense Promoter of L1 Retrotransposon. *Genomics*, **79**, 628–634.
- 38 15. Miglio, U., Berrino, E., Panero, M., Ferrero, G., Coscujuela Tarrero, L., Miano, V., Dell’Aglío,  
39 C., Sarotto, I., Annaratone, L., Marchiò, C., *et al.* (2018) The expression of LINE1-MET chimeric  
40 transcript identifies a subgroup of aggressive breast cancers. *Int. J. Cancer*, **143**, 2838–2848.
- 41 16. Hur, K., Cejas, P., Feliu, J., Moreno-Rubio, J., Burgos, E., Boland, C.R. and Goel, A. (2014)  
42 Hypomethylation of long interspersed nuclear element-1 (LINE-1) leads to activation of proto-  
43 oncogenes in human colorectal cancer metastasis. *Gut*, **63**, 635–646.
- 44 17. Wolff, E.M., Byun, H.-M., Han, H.F., Sharma, S., Nichols, P.W., Siegmund, K.D., Yang, A.S.,  
45 Jones, P.A. and Liang, G. (2010) Hypomethylation of a LINE-1 promoter activates an alternate  
46 transcript of the MET oncogene in bladders with cancer. *PLoS Genet.*, **6**, e1000917.
- 47 18. Mätlik, K., Redik, K. and Speek, M. (2006) L1 Antisense Promoter Drives Tissue-Specific  
48 Transcription of Human Genes. *J. Biomed. Biotechnol.*, **2006**, 1–16.
- 49 19. Cruickshanks, H.A., Vafadar-Isfahani, N., Dunican, D.S., Lee, A., Sproul, D., Lund, J.N.,  
50 Meehan, R.R. and Tufarelli, C. (2013) Expression of a large LINE-1-driven antisense RNA is  
51 linked to epigenetic silencing of the metastasis suppressor gene TFPI-2 in cancer. *Nucleic Acids*  
52 *Res.*, **41**, 6857–6869.
- 53 20. Cruickshanks, H.A. and Tufarelli, C. (2009) Isolation of cancer-specific chimeric transcripts  
54 induced by hypomethylation of the LINE-1 antisense promoter. *Genomics*, **94**, 397–406.

- 1 21. Babaian, A., Thompson, I.R., Lever, J., Gagnier, L., Karimi, M.M. and Mager, D.L. (2019)  
2 LIONS: analysis suite for detecting and quantifying transposable element initiated transcription  
3 from RNA-seq. *Bioinforma. Oxf. Engl.*, **35**, 3839–3841.
- 4 22. Jang, H.S., Shah, N.M., Du, A.Y., Dailey, Z.Z., Pehrsson, E.C., Godoy, P.M., Zhang, D., Li, D.,  
5 Xing, X., Kim, S., *et al.* (2019) Transposable elements drive widespread expression of oncogenes  
6 in human cancers. *Nat. Genet.*, **51**, 611–617.
- 7 23. Pinson, M.-E., Pogorelcnik, R., Court, F., Arnaud, P. and Vaurs-Barrière, C. (2018) CLIFinder:  
8 identification of LINE-1 chimeric transcripts in RNA-seq data. *Bioinformatics*, **34**, 688–690.
- 9 24. Louis, D.N., Perry, A., Reifenberger, G., von Deimling, A., Figarella-Branger, D., Cavenee, W.K.,  
10 Ohgaki, H., Wiestler, O.D., Kleihues, P. and Ellison, D.W. (2016) The 2016 World Health  
11 Organization Classification of Tumors of the Central Nervous System: a summary. *Acta*  
12 *Neuropathol. (Berl.)*, **131**, 803–820.
- 13 25. Cohen, A.L., Holmen, S.L. and Colman, H. (2013) IDH1 and IDH2 mutations in gliomas. *Curr.*  
14 *Neurol. Neurosci. Rep.*, **13**, 345.
- 15 26. Chen, J., Li, Y., Yu, T.-S., McKay, R.M., Burns, D.K., Kernie, S.G. and Parada, L.F. (2012) A  
16 restricted cell population propagates glioblastoma growth after chemotherapy. *Nature*, **488**, 522–  
17 526.
- 18 27. Lathia, J.D., Mack, S.C., Mulkearns-Hubert, E.E., Valentim, C.L.L. and Rich, J.N. (2015) Cancer  
19 stem cells in glioblastoma. *Genes Dev.*, **29**, 1203–1217.
- 20 28. Hedges, L.V. (1981) Distribution Theory for Glass’s Estimator of Effect Size and Related  
21 Estimators. *J. Educ. Stat.*, **6**, 107.
- 22 29. Oakes, C.C., La Salle, S., Robaire, B. and Trasler, J.M. (2006) Evaluation of a Quantitative DNA  
23 Methylation Analysis Technique using Methylation-Sensitive/Dependent Restriction Enzymes  
24 and Real-Time PCR. *Epigenetics*, **1**, 146–152.
- 25 30. Lavie, L. (2004) The human L1 promoter: Variable transcription initiation sites and a major  
26 impact of upstream flanking sequence on promoter activity. *Genome Res.*, **14**, 2253–2260.
- 27 31. Philippe, C., Vargas-Landin, D.B., Doucet, A.J., van Essen, D., Vera-Otarola, J., Kuciak, M.,  
28 Corbin, A., Nigumann, P. and Cristofari, G. (2016) Activation of individual L1 retrotransposon  
29 instances is restricted to cell-type dependent permissive loci. *eLife*, **5**, e13926.
- 30 32. Lock, F.E., Babaian, A., Zhang, Y., Gagnier, L., Kuah, S., Weberling, A., Karimi, M.M. and  
31 Mager, D.L. (2017) A novel isoform of IL-33 revealed by screening for transposable element  
32 promoted genes in human colorectal cancer. *PLoS One*, **12**, e0180659.
- 33 33. Vafadar-Isfahani, N., Parr, C., McMillan, L.E., Sanner, J., Yeo, Z., Saddington, S., Peacock, O.,  
34 Cruickshanks, H.A., Meehan, R.R., Lund, J.N., *et al.* (2017) Decoupling of DNA methylation and  
35 activity of intergenic LINE-1 promoters in colorectal cancer. *Epigenetics*, **12**, 465–475.
- 36 34. Court, F., Le Boiteux, E., Fogli, A., Müller-Barthélémy, M., Vaurs-Barrière, C., Chautard, E.,  
37 Pereira, B., Biau, J., Kemeny, J.-L., Khalil, T., *et al.* (2019) Transcriptional alterations in glioma  
38 result primarily from DNA methylation-independent mechanisms. *Genome Res.*, **29**, 1605–1621.
- 39 35. Le Boiteux, E., Court, F., Guichet, P., Vaurs-Barrière, C., Vaillant, I., Chautard, E., Verrelle, P.,  
40 Costa, B.M., Karayan-Tapon, L., Fogli, A., *et al.* (2021) Widespread overexpression from the  
41 four DNA hypermethylated HOX clusters in aggressive (*IDH* wt) glioma is associated with  
42 H3K27me3 depletion and alternative promoter usage. *Mol. Oncol.*, **15**, 1995–2010.
- 43 36. Cohen, A.L., Holmen, S.L. and Colman, H. (2013) IDH1 and IDH2 Mutations in Gliomas. *Curr.*  
44 *Neurol. Neurosci. Rep.*, **13**, 345.
- 45 37. Villalva, C., Martin-Lannerée, S., Cortes, U., Dkhissi, F., Wager, M., Le Corf, A., Tourani, J.-M.,  
46 Dusanter-Fourt, I., Turhan, A.G. and Karayan-Tapon, L. (2011) STAT3 is essential for the  
47 maintenance of neurosphere-initiating tumor cells in patients with glioblastomas: A potential for  
48 targeted therapy? *Int. J. Cancer*, **128**, 826–838.
- 49 38. Villalva, C., Cortes, U., Wager, M., Tourani, J.-M., Rivet, P., Marquant, C., Martin, S., Turhan,  
50 A.G. and Karayan-Tapon, L. (2012) O6-Methylguanine-Methyltransferase (MGMT) Promoter  
51 Methylation Status in Glioma Stem-Like Cells is Correlated to Temozolomide Sensitivity Under  
52 Differentiation-Promoting Conditions. *Int. J. Mol. Sci.*, **13**, 6983–6994.
- 53 39. Guichet, P.-O., Masliantsev, K., Tachon, G., Petropoulos, C., Godet, J., Larrieu, D., Milin, S.,  
54 Wager, M. and Karayan-Tapon, L. (2018) Fatal correlation between YAP1 expression and glioma  
55 aggressiveness: clinical and molecular evidence: YAP1 expression in gliomas. *J. Pathol.*, **246**,  
56 205–216.
- 57



## 1 **Legends to Figures**

2

### 3 **Figure 1: LCTs are produced from the recent L1PA1 to L1PA7 L1 subfamilies**

4 **a.** Percentage of L1 elements associated with CLIFinder-detected chimeras (red) among all L1  
5 elements with a 5' UTR region and belonging to the human L1PA1 to L1PA8 subfamilies. **b.** RT-  
6 PCR-based 5' L1 walking approach. This technique combines a LCT-specific primer (LCT-US, green  
7 bar), located in the unique genomic sequence (in green), and primers (grey lines) located either  
8 downstream or upstream of the L1-ASP 200 and 450bp regions (highlighted boxes) in the L1 5'UTR  
9 sequence (in blue). **c.** Examples of RT-PCR 5' L1 walking results in two glioma samples (GS#1 and  
10 #2). The chimeras Id\_1087 (L1PA2 subfamily) and Id\_179 (L1PA7 subfamily) correspond to LCTs  
11 with a transcription start site (TSS) in the L1-ASP-450 region. The chimera Id\_2900 (L1PA8  
12 subfamily) initiates upstream the two L1-ASP regions and therefore, it is not considered to be an LCT.  
13 **d.** Results obtained for the 36 tested chimeras distributed according to the L1 subfamilies. The TSS of  
14 each chimera is defined according to the PCR results obtained with the L1 downstream and upstream  
15 primers. Chimeras are defined as LCTs when no amplification was obtained with the L1 upstream  
16 primer on the tumor cDNA samples. **e.** The long-read SMRT sequencing dataset was used to generate  
17 a heatmap showing TSS occurrence along the full-length L1 element associated to the 226 putative  
18 LCTs defined by CLIFinder. Twenty-two LCTs were not detected in the long-read SMRT sequencing  
19 dataset and are shown as white horizontal lines. Among the other 204 LCTs, 84.80% had at least one  
20 read with a TSS included in the 0-500bp 5'UTR region, and this percentage increased to 98% when  
21 considering TSS in the whole 5'UTR.

22

### 23 **Figure 2: LCTs can be polyA and non-polyA transcripts**

24 **a.** Heatmap describing the PCR amplification results for 24 LCTs using glioma cDNA obtained with  
25 different reverse transcription conditions (see Methods). PCR amplifications were performed using a  
26 LCT-specific primer (LCT-US) and the L1 common primer (see Fig. 1B). **b.** Comparison of the  
27 number of LCTs (identified by CLIFinder) for each sample (glioma samples and Glioma Stem Cell

1 (GSC) lines) and according to the RNA-seq approach used (Ribo-Zero (RBZ) and PolyA) \*: This  
2 GSC2 cell line was analyzed by PolyA- and RBZ-RNA-seq. **c.** Comparison of the number of LCTs  
3 identified by CLIFinder in three datasets: PolyA RNA-seq dataset of two GSC cell lines (GSC2 and  
4 GSC6), Ribo-Zero RNA-seq dataset of two GSC lines (GSC2 and GSC1), and Ribo-Zero RNA-seq  
5 dataset of eight *IDHwt* glioma samples. **d.** Comparison of the number of LCTs (identified by  
6 CLIFinder) associated with L1PA1 to L1PA7 retrotransposons in the RBZ and in PolyA RNA-seq  
7 datasets of the GSC2 line (data obtained using the same RNA sample).

8

9 **Figure 3: Characteristics of the 1,509 LCTs identified by CLIFinder in *IDHwt* glioma and**  
10 **normal brain samples**

11 **a.** Comparison of the localization, relative to annotated genes, of all recent L1 with a 5' UTR  
12 (n=8,744) and of recent L1 retrotransposons associated with the 1,509 LCTs. The L1-ASP orientation  
13 is given relative to the gene orientation. The significant enrichment and depletion of intragenic and  
14 intergenic LCT-producing L1 retrotransposons, respectively, compared with all recent L1  
15 retrotransposons, was confirmed with the binomial test (\*\*p <0.001). **b.** Number of genes in which  
16 the indicated number of LCTs was identified by CLIFinder; \*\*p <0.005 (one way-ANOVA and  
17 Mann-Whitney post hoc test). **c.** Gene Ontology terms (GOTERM) enriched in genes containing  
18 recent L1 retrotransposons associated with LCTs (LCT+, red) compared with genes containing recent  
19 L1 retrotransposons that do not produce LCTs (LCT-, black). For each GOTERM category, the most  
20 significant terms (Bonferroni corrected p <0.05) are shown.

21

22 **Figure 4: LCTs are expressed in normal brain and glioma samples, but with different expression**  
23 **levels for a LCT subset.**

24 **a.** Heatmap showing the centered-scaled expression levels between samples for each LCT and the  
25 differential expression of the 31 validated LCTs in *IDHwt* glioma and control samples (Mann-  
26 Whitney, MW, test followed by the Bonferroni-Holm, BH, correction). Significant LCT expression  
27 up- or down-regulation in glioma samples is represented by red and green squares, respectively. **b.**

1 LCT\_718 (shown by an arrow in a) is an example of LCT significantly overexpressed in glioma  
2 samples. (\*\*\*) BH adjusted  $p < 0.005$ ) **c.** Positive Spearman's correlation between the CLIFinder reads  
3 number for the 31 validated LCTs and their relative expression ( $\log_{10}$ ) measured by RT-qPCR in two  
4 control and eight *IDHwt* glioma samples. **d.** Hedges' Effect Size (ES) determination. For each of the  
5 31 validated LCTs, the Hedges's ES (g value) between *IDHwt* glioma and control samples was  
6 calculated for the CLIFinder reads number and RT-qPCR relative expression. An ES cut-off value of 1  
7 for the CLIFinder reads number was determined with a sensitivity of 44% and a specificity of 85%. **e.**  
8 Heatmap representations of the centered-scaled LCT CLIFinder reads number and of the predicted  
9 differential expression in *IDHwt* glioma *versus* control samples for the 1,509 LCTs. Hedges' g values  
10  $< -1$  and  $> +1$  predict a significant up- (red) and down-regulation (green), respectively. Values between  
11 these cut-offs predict no expression change (beige). For each expression group (up- and down-  
12 regulation), the localization of the LCT-producing L1 element, relative to the annotated genes, is  
13 shown in the pie charts: intragenic LCTs transcribed in the sense of the gene (Gene-S) in black,  
14 intragenic LCTs transcribed in the opposite sense (Gene\_AS) in grey, and intergenic LCTs in white.

15

16 **Figure 5: L1 hypomethylation and CNVs are not the driving force of LCT expression changes in**  
17 ***IDHwt* glioma samples**

18 **a.** Heatmaps summarizing the expression and methylation index of the L1 5'UTR that produce the six  
19 indicated LCTs. The correlation analysis between these parameters is shown on the right panel. **b.**  
20 Detailed results obtained for LCT\_887 (black arrow in a), in control and glioma samples. Although the  
21 L1 promoter methylation index is significantly lower in *IDHwt* glioma than in control samples, the  
22 Spearman's correlation analysis did not highlight any negative correlation between L1 promoter  
23 methylation index and LCT\_887 expression level. **c.** Heatmaps summarizing LCT expression and  
24 CNV data for the 31 validated LCT loci. The results of the Spearman's correlation analysis are shown  
25 in the right panel. **d.** Example of results obtained for the LCTs\_1425 and\_1440-1441. These two LCTs  
26 are on chromosome 7 that is duplicated in *IDHwt* glioma; however, the expression level of LCT\_1425  
27 was not changed in glioma compared with control samples, and no correlation was detected between  
28 LCT\_1440-1441 expression level and CNV.

1

2 **Figure 6: Recent LCT-producing L1 retrotransposons are localized in transcriptionally active**  
3 **regions**

4 **a.** Comparison of genes containing a recent L1 retrotransposon associated or not with the transcription  
5 of a LCT in this study. 270 genes contain at least two recent L1 retrotransposons with divergent LCT  
6 expression status. **b.** Comparison of the expression levels of genes containing a recent L1  
7 retrotransposon that expresses (Genes L1+LCT+, n=433) or not (Genes L1+LCT-, n=969) LCTs. The  
8 box plots lower and upper limits indicate the 25<sup>th</sup> and 75<sup>th</sup> percentile, respectively, and the middle line  
9 represents the median. The median fpkm numeric value is given above each group; \*\*\*p <2.2e-16 (T-  
10 test). **c.** Analysis of ChIP-seq reads density data for input, H3K36me3 and H3K9me3 for two normal  
11 brain samples and one GSC cell line. Plots are centered on a ±1 kb window according to the L1 start  
12 and end (blue horizontal bar with the L1-ASP in red). Both intragenic (dark green) and intergenic  
13 (light green) L1 retrotransposons are considered. Loci retained as recent L1 retrotransposons not  
14 associated with LCT transcription (dashed lines) correspond to L1 elements never associated with  
15 LCT in our study (n=1,332 intragenic and n= 5,359 intergenic L1). For analysis of normal brain  
16 samples, all recent L1 retrotransposons associated with CLIFinder-detected LCTs in controls and  
17 *IDHwt* gliomas were retained (n=424 intragenic L1 and n=635 intergenic L1 retrotransposons). For  
18 analysis of the GSC sample, only recent L1 retrotransposons associated with LCTs identified by  
19 CLIFinder in GSC samples after Ribo-Zero RNA-seq were analyzed (n=232 intragenic L1 and n=58  
20 intergenic L1 retrotransposons). Plots represent the mean ChIP-seq signal values for recent L1  
21 retrotransposons associated with LCT transcription (solid lines) and for recent L1 retrotransposons not  
22 associated with LCT transcription (dashed lines).

23

24 **Figure 7: LCT deregulation in glioma samples is linked to transcriptional changes in the host**  
25 **gene**

26 **a.** Heatmaps summarizing the expression levels quantified by RT-qPCR of the indicated LCTs (13  
27 intragenic LCTs, including 11 overexpressed in *IDHwt* glioma samples) and of their host genes. The

1 host gene expression levels were compared between Controls and *IDHwt* glioma samples and  
2 correlated with the LCT relative expression level. Significance of Spearman correlations were assessed  
3 by Bonferroni corrected p-values and indicated as \*p <0.05, \*\*p <0.01. **b.** Example of results obtained  
4 for LCT\_886 and its host gene *CCDC109B*. LCT\_886 is overexpressed in *IDHwt* glioma samples.  
5 *CCDC109B* also is significantly overexpressed in *IDHwt* glioma samples compared with controls, and  
6 LCT\_886 and *CCDC109B* expression are positively correlated (Spearman's rho = 0.814; p-value = 0).  
7 **c.** Comparison of the expression level (fpkm) in control (n=3) and *IDHwt* glioma samples (n=8) of the  
8 genes associated with the 64 and 120 intragenic LCTs predicted to be up- or down-regulated,  
9 respectively, in tumors by the Hedges' ES prediction model (Fig. 4E).

10

11 **Figure 8: Working model of LCT transcriptional regulation in normal brain and in**  
12 **glioma/GSCs**

13 Our study supports a model whereby the transcriptional activity of the host genomic locus is the main  
14 driving force of LCT expression. We propose that this model is applicable particularly when the host  
15 gene and the L1-ASP are in the same orientation. Accordingly, in normal brain, LCTs are widely  
16 expressed. Changes in the host gene/genomic locus expression level in glioma and GSCs will lead to  
17 similar changes in LCT expression level, regardless of the L1 promoter methylation status. Note that  
18 according to this model, a gene specifically expressed in glioma can contribute to the expression of a  
19 glioma-specific LCT (locus E).

20

21

1 **Abbreviations**

2 ASP: antisense promoter

3 CLIFinder: Chimeric Line Finder

4 CNV: copy number variation

5 ES: Effect size

6 GOTERM: Gene Ontology Terms

7 GS: Glioma Sample

8 GSC: glioma stem cells

9 IDHwt: wild type Isocitrate Dehydrogenase gene

10 L1: long interspersed element 1

11 LCT: L1 chimeric transcripts

12 LCT-US: LCT-unique sequence

13 RBZ: Ribo-Zero RNA-seq

14 SMRT: single molecule real time sequencing

15 TSS: Transcription start site

16