

Supplemental Material

Uncovering transposable element variants and their potential adaptive impact in urban populations of the malaria vector *Anopheles coluzzii*

Carlos Vargas-Chavez¹, Neil Michel Longo Pandy^{2,3}, Sandrine E. Nsango⁴, Laura Aguilera¹, Diego Ayala^{2,5*}, Josefa González^{1*}

¹Institute of Evolutionary Biology (CSIC-Universitat Pompeu Fabra), Barcelona, Spain.

²CIRMF, Franceville, Gabon.

³Ecole doctorale en infectiologie tropicale (EDR), Franceville, Gabon.

⁴ Faculté de Médecine et des Sciences Pharmaceutiques, Université de Douala, BP 2701, Douala, Cameroun

⁵MIVEGEC, Univ Montpellier, CNRS, IRD, Montpellier, France

Emails: carlos.vargas@ibe.upf-csic.es, longo2michel@gmail.com, nsango2013@yahoo.fr, mlaura.aguilera@gmail.com, diego.ayala@ird.fr, josefa.gonzalez@csic.es

*Corresponding authors

Contents

Supplemental Figures	2
Figure S1. Number of TE copies identified when using the TE libraries from an increasing number of genomes.....	2
Figure S2. Novel TE families.....	3
Figure S3. Number of TE insertions vs. genome size.....	5
Figure S4. Differences in TE content between the seven <i>An. coluzzii</i> genomes.....	6
Figure S5. TE landscapes.....	8
Figure S6. Genes with TE insertions from active families.....	9
Figure S7. Genes associated with insecticide resistance with TE insertions.....	74
Figure S8. TE abundance plots by chromosome.....	95
Supplemental Table Descriptions	97
Supplemental Files	100
Supplemental Methods	101
Supplemental References	107

Supplemental Figures

Figure S1. Number of TE copies identified when using the TE libraries from an increasing number of genomes. Number of TE copies identified when using the TE library of a single genome or when using all possible combinations of more than one genome.

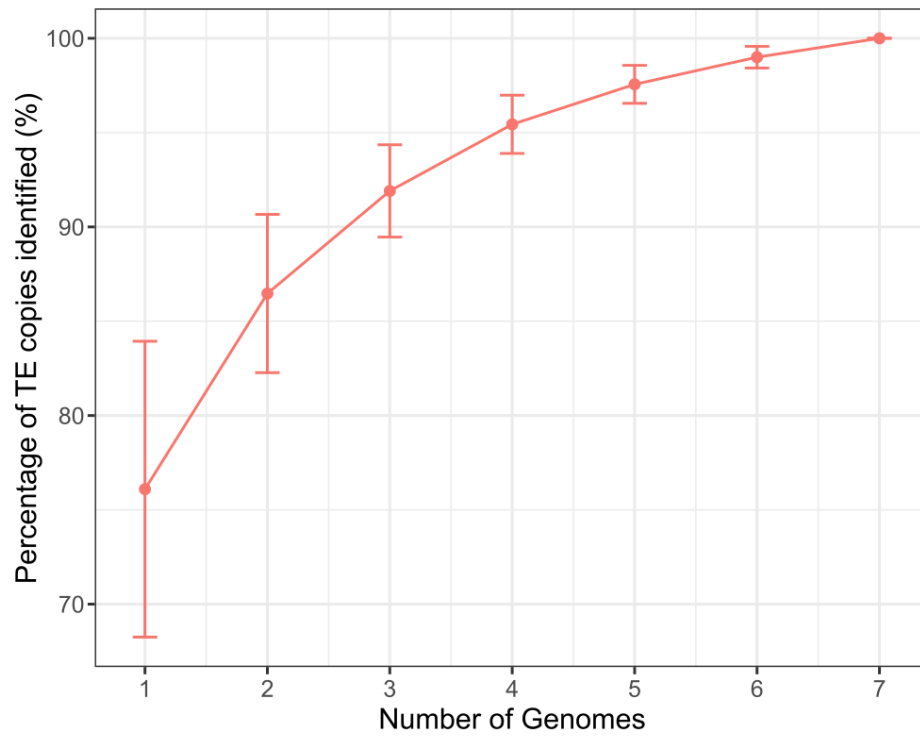
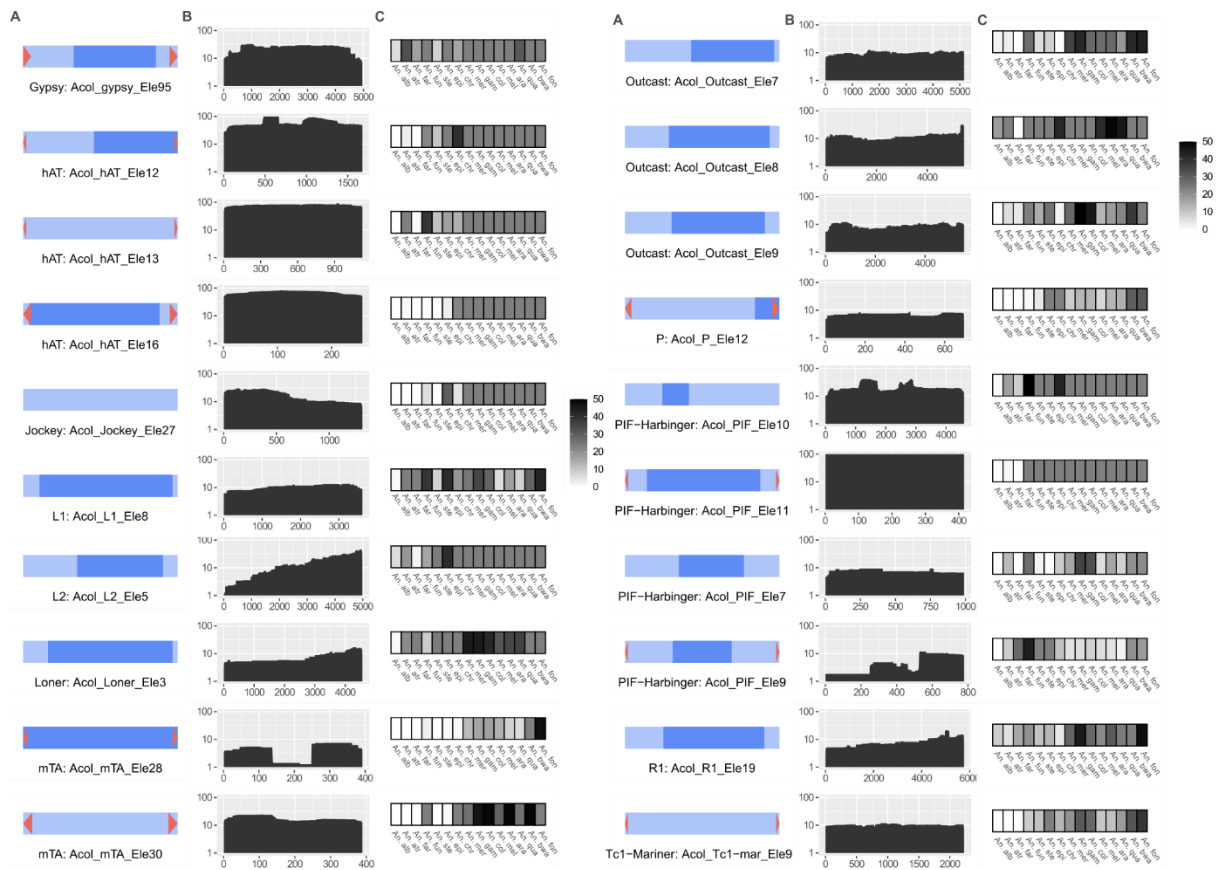


Figure S2. Novel TE families. A) The structure of each new family is displayed: the light blue box represents the full extension of the TE and the red arrows represent LTRs. B) All insertions for each TE family are shown as a coverage plot where each line represents a copy in a genome. C) Phylogenetic distribution of the TE family insertions in 15 members of the *Anopheles* genus, *Culex quinquefasciatus*, *Aede aegypti* and *D. melanogaster*. The number of insertions with more than 80% identity and spanning at least 80% of the consensus sequence, in each species is shown using a black and white gradient: species with no insertions are shown in white while species with 50 or more insertions are shown in black.



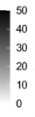
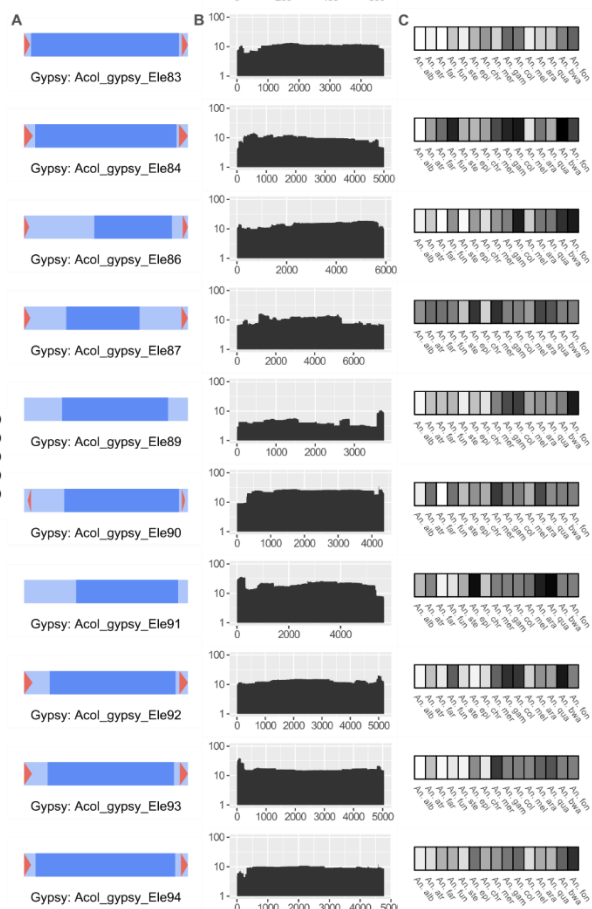
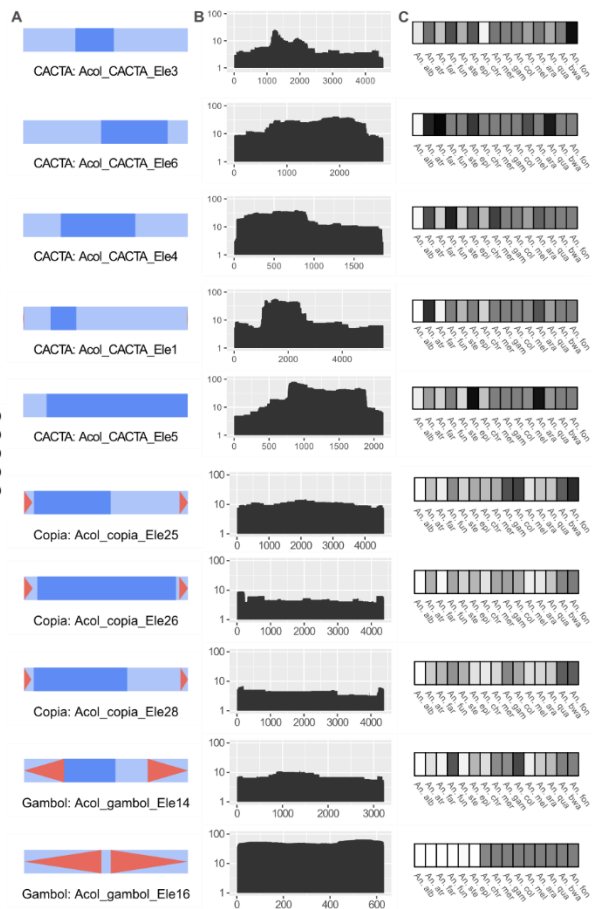
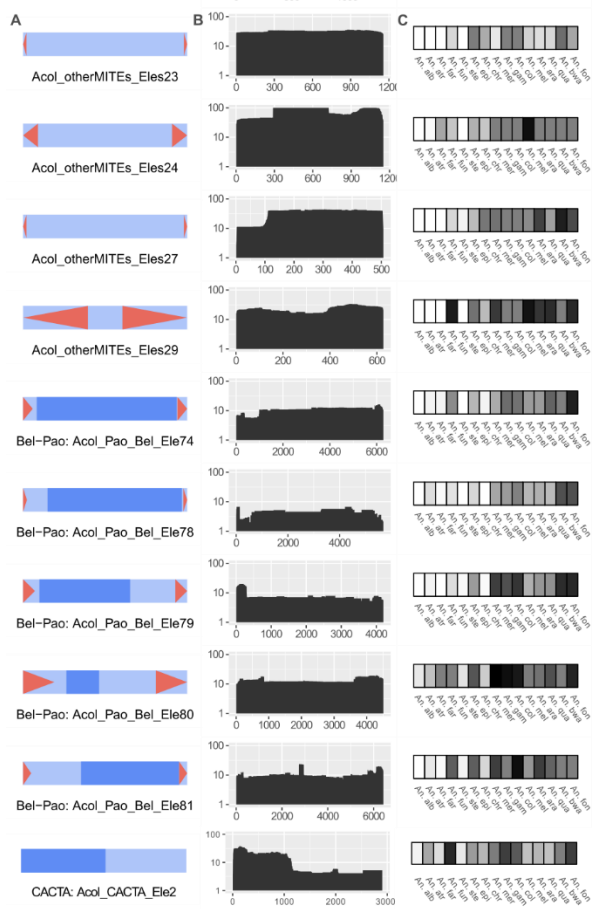
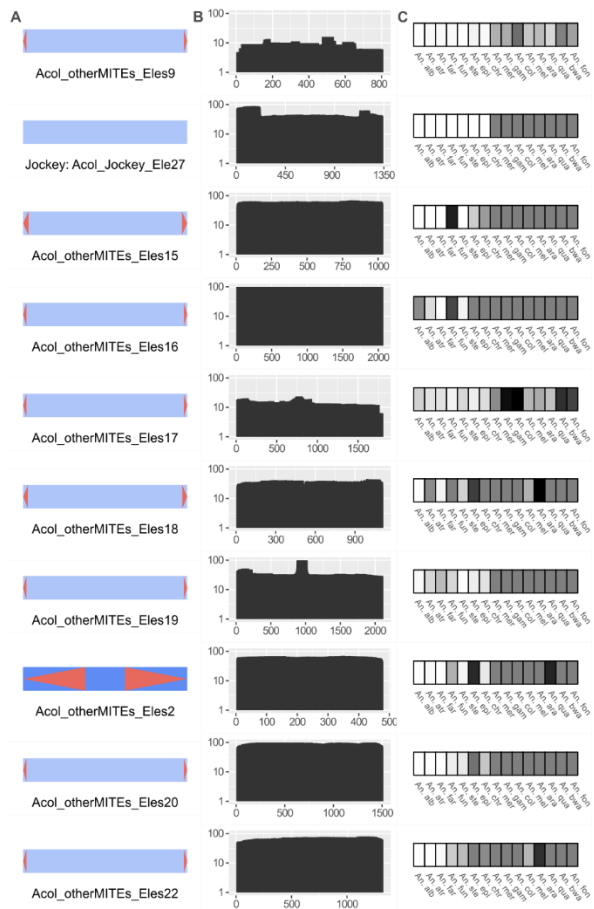


Figure S3. Number of TE insertions vs. genome size. Comparison of the bases spanned by TEs in each genome with their full genome sizes.

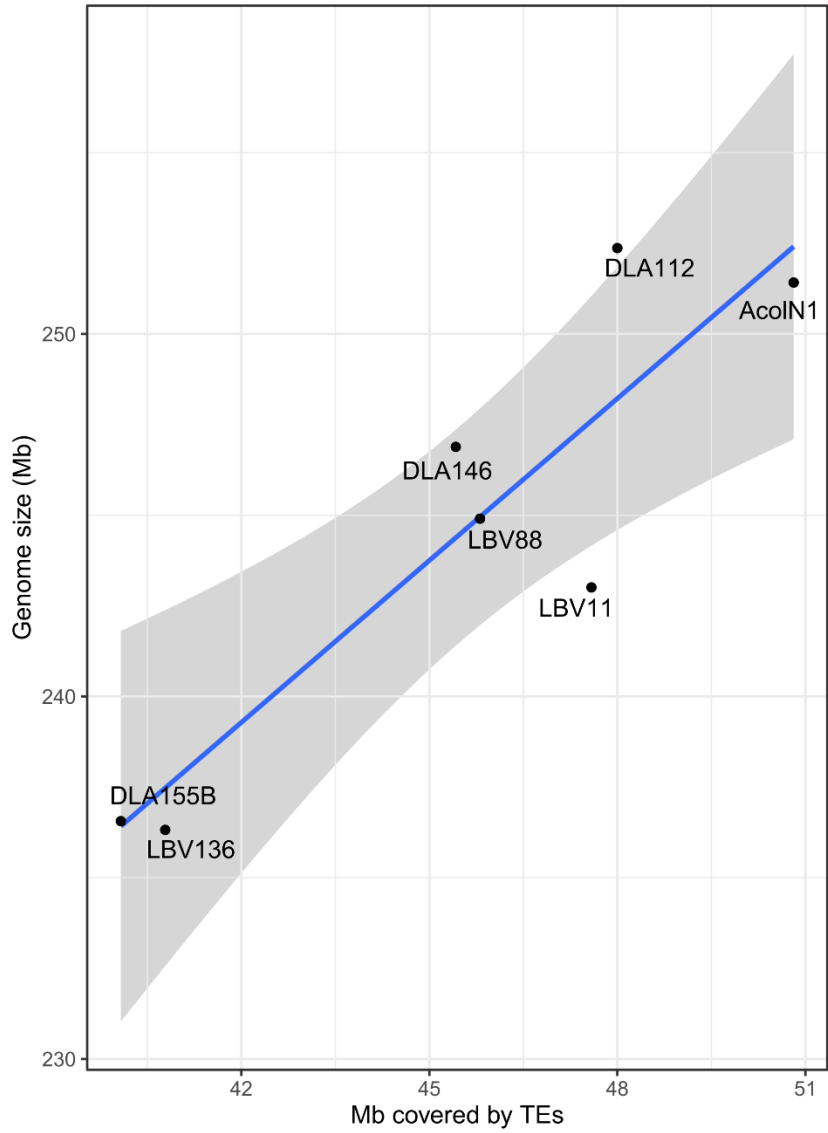
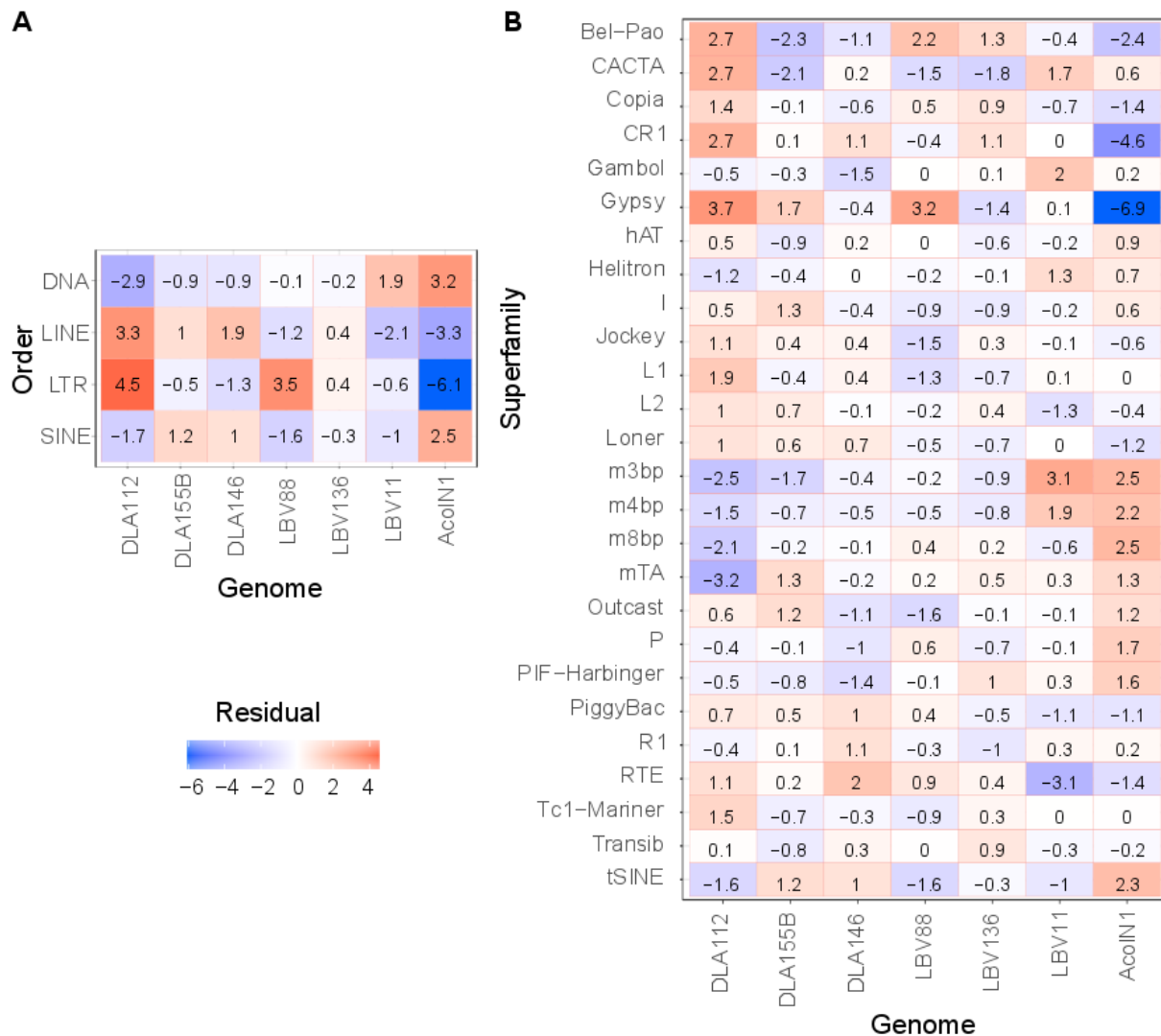


Figure S4. Differences in TE content between the seven *An. coluzzii* genomes.

Differences are shown at the (A) order and (B) superfamily levels. χ^2 tests were performed for the number of insertions and the Pearson's residuals are shown. Note that MITEs are divided into the *m3bp*, *m4bp*, *m8bp* and *mTA* superfamilies. (C) PCA using the SNPs identified using the Illumina data selecting (i) those that were found on introns (and thus most likely to be neutral and to reflect the population structure) and (ii) those located on chromosome 3, as this chromosome is free from polymorphic inversions and regions of reduced recombination (Miles et al 2017, Science)



C

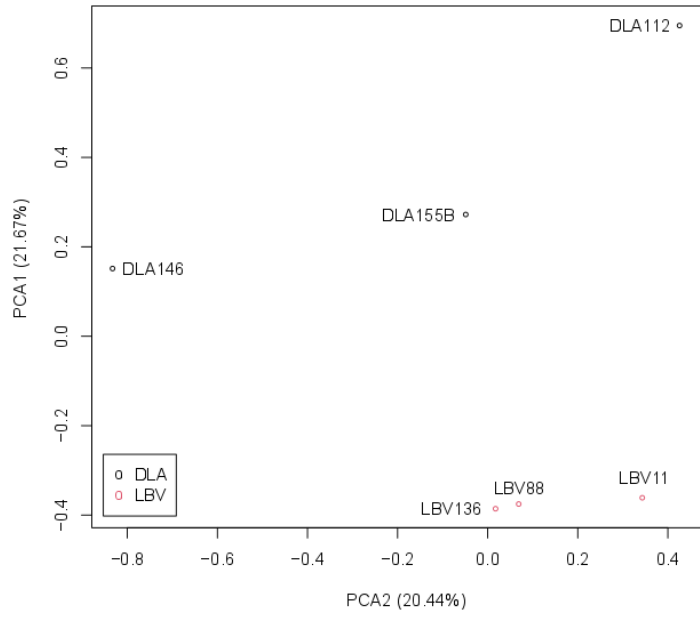


Figure S5. TE landscapes. TE landscapes for the six genomes sequenced in this work generated using dnaPipeTE.

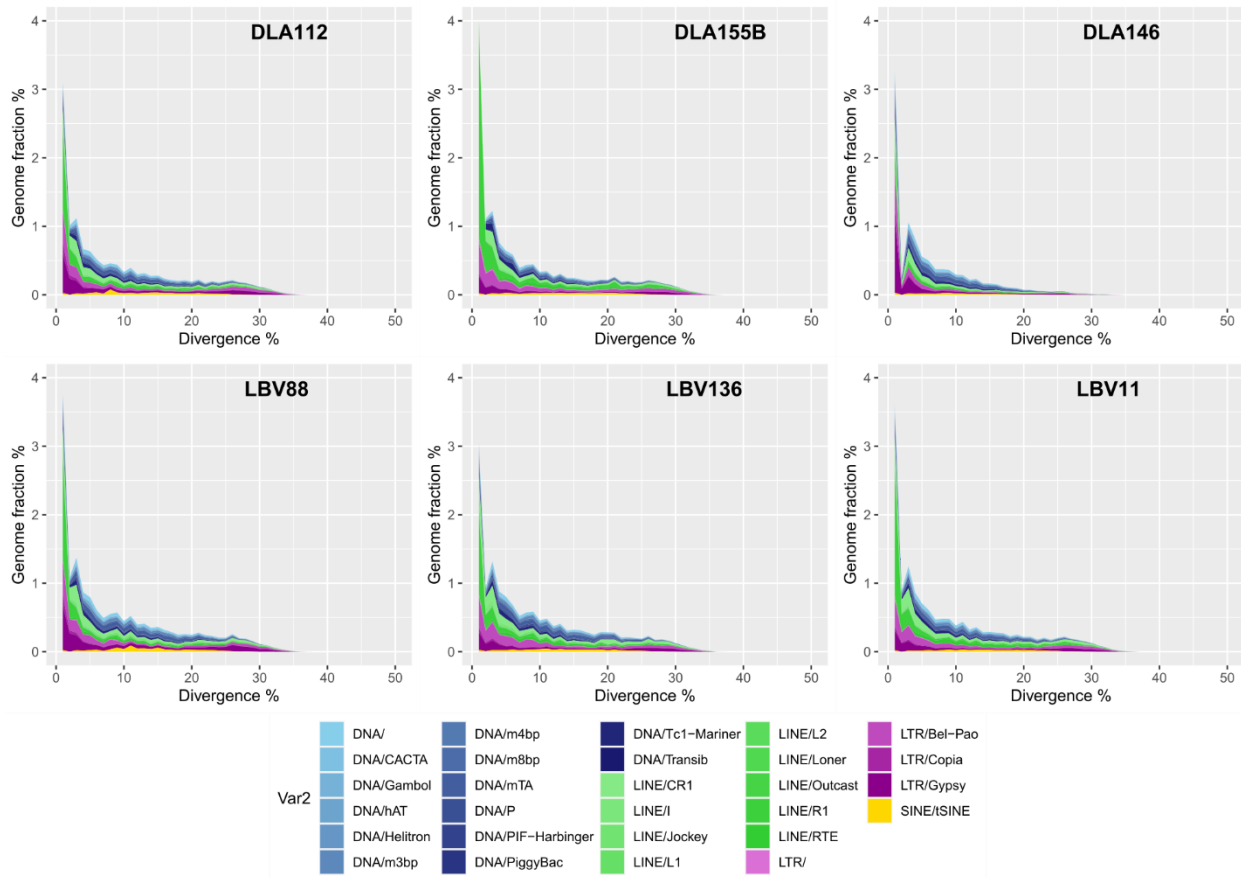
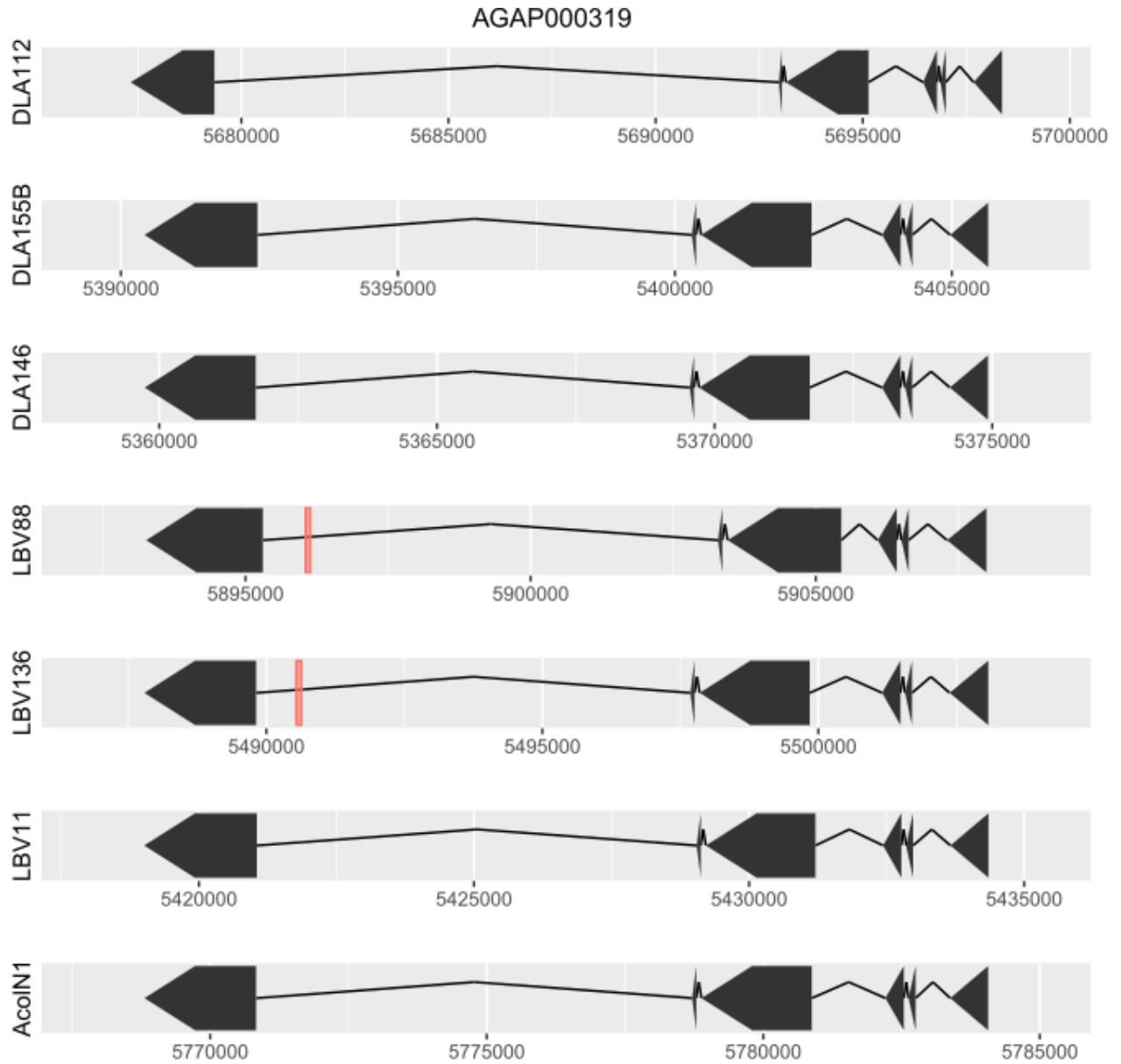
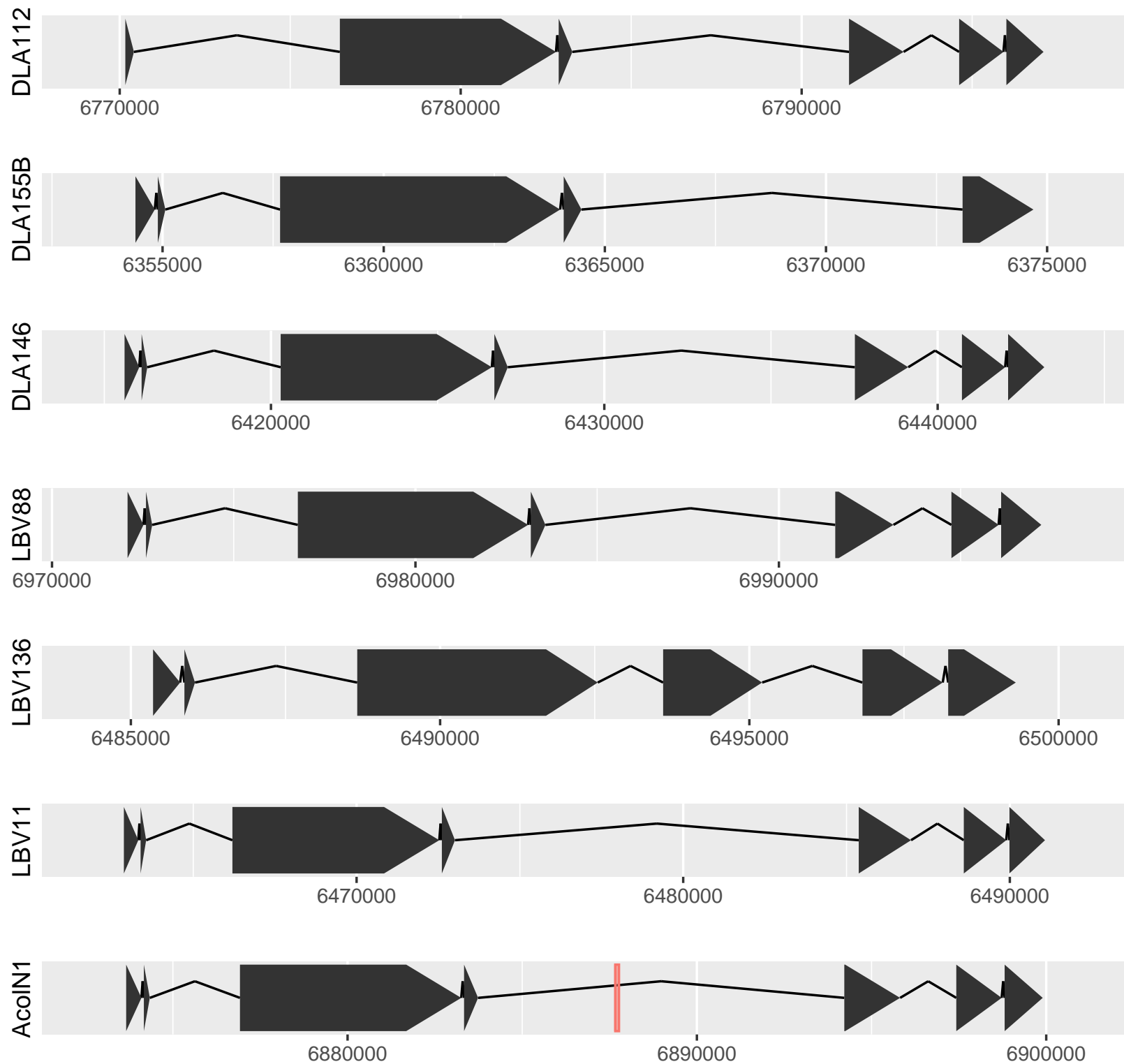


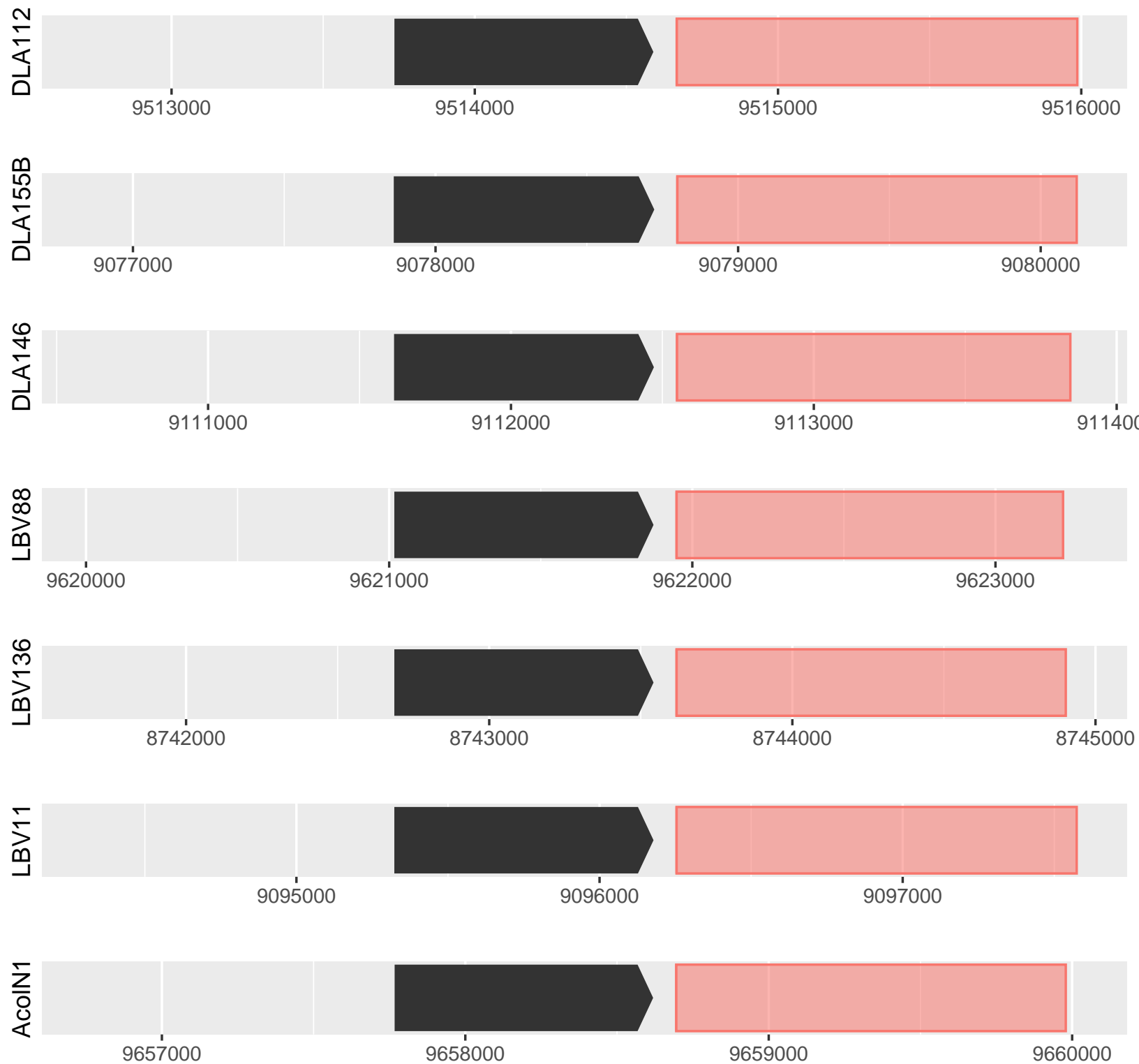
Figure S6. Genes with TE insertions from active families. Diagrams of TE insertions located inside genes or less than 1 kb from a gene showing the gene structure and the TE insertion (in red).



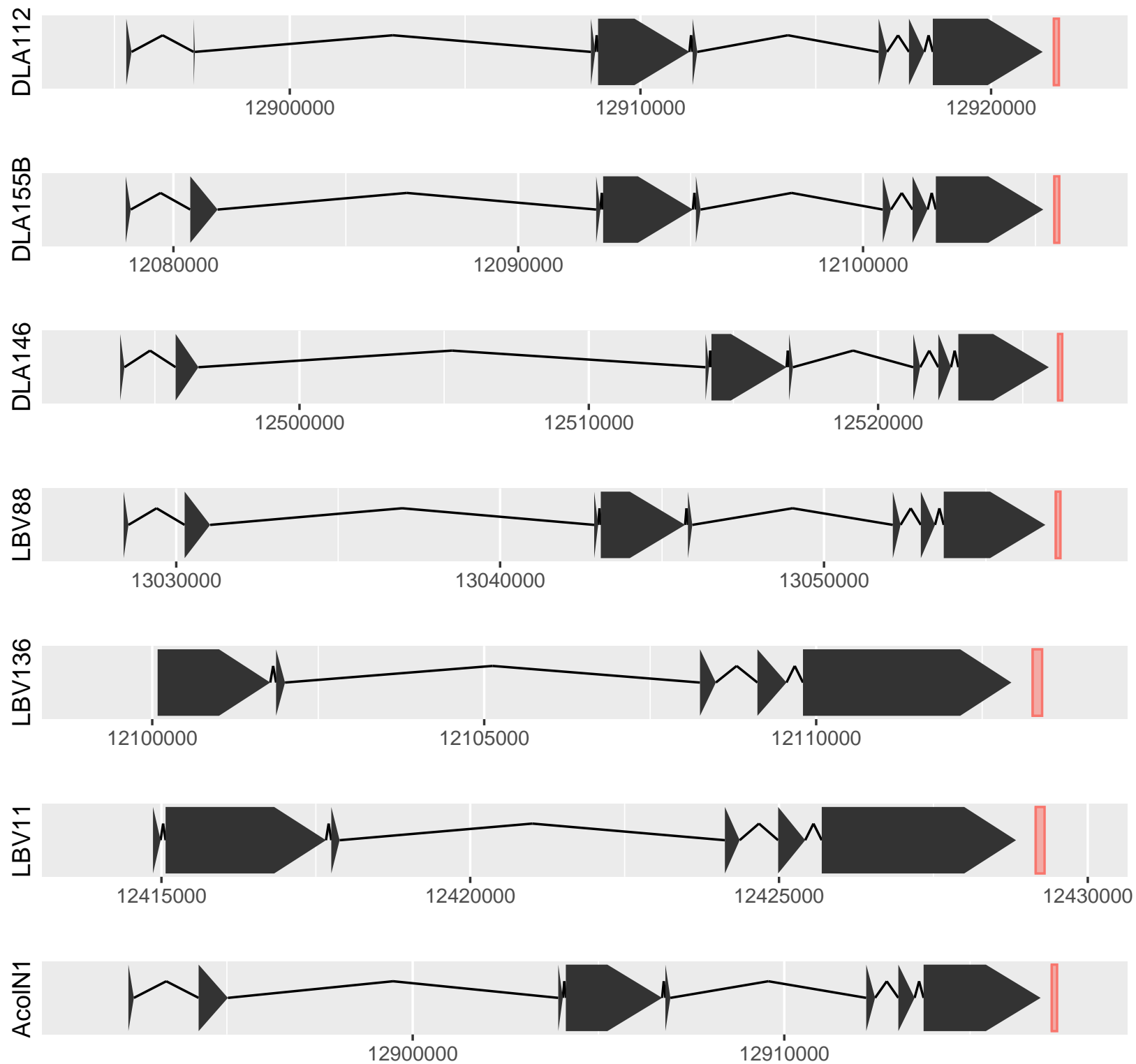
AGAP000368



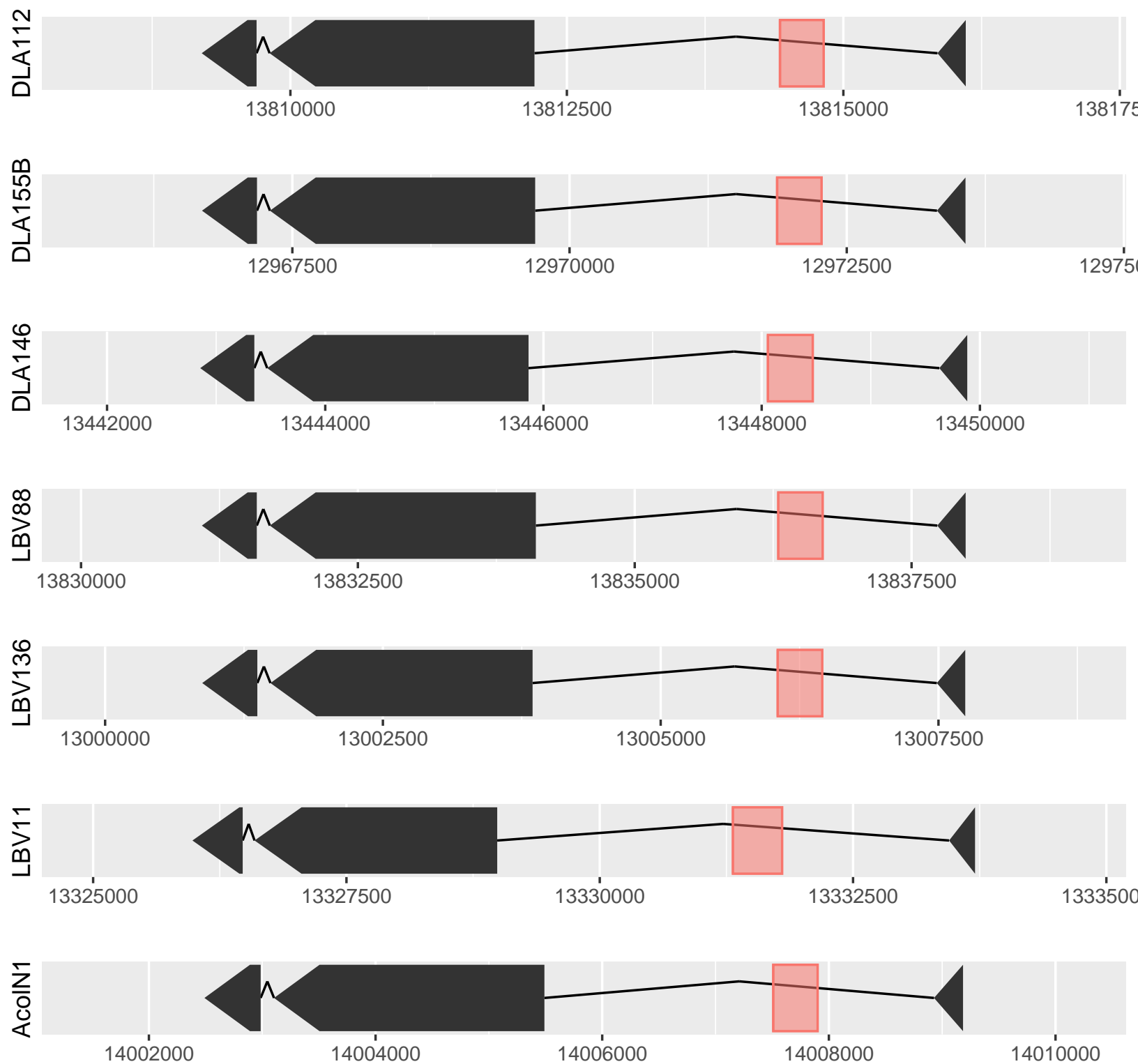
AGAP000531



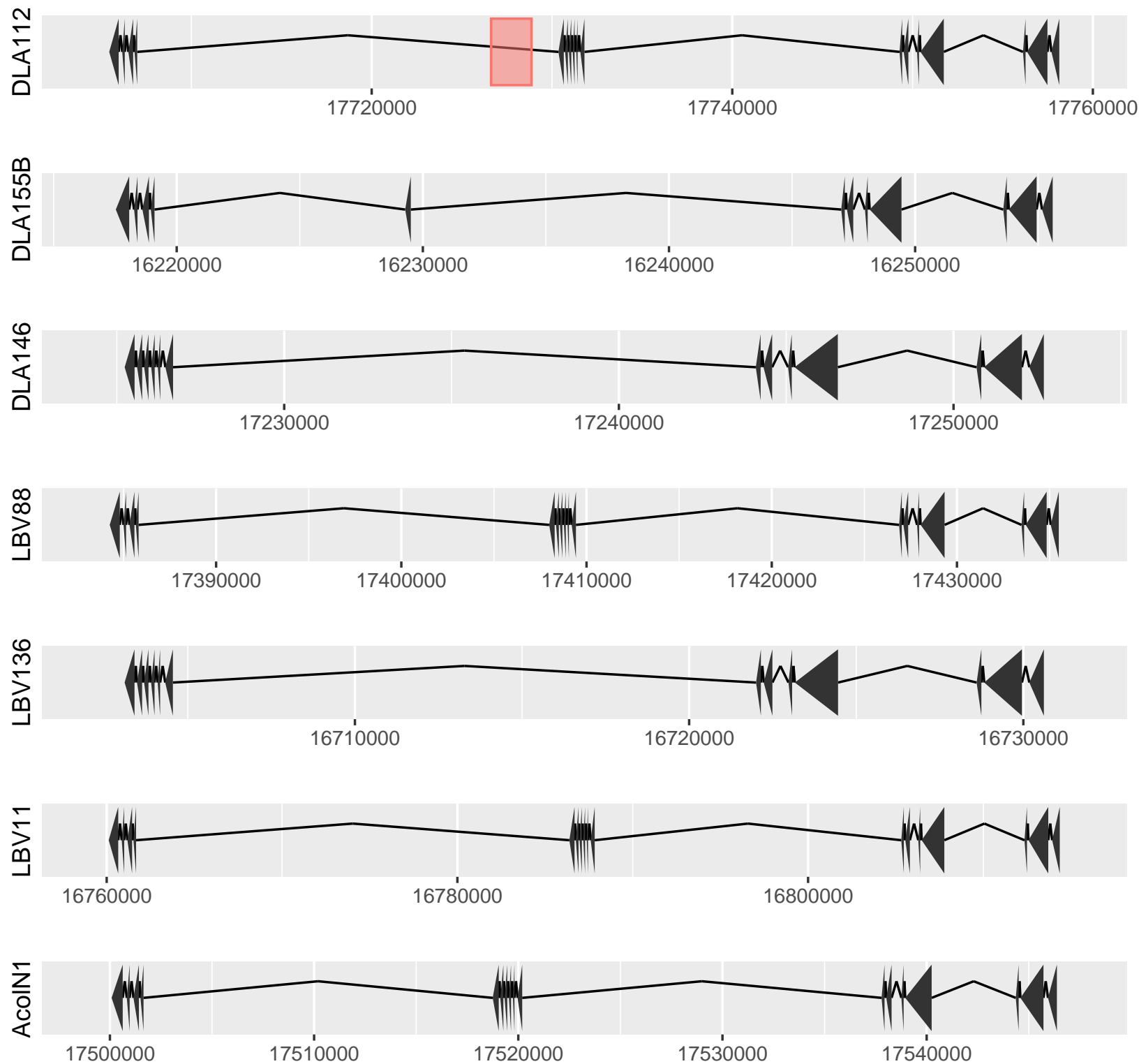
AGAP000707



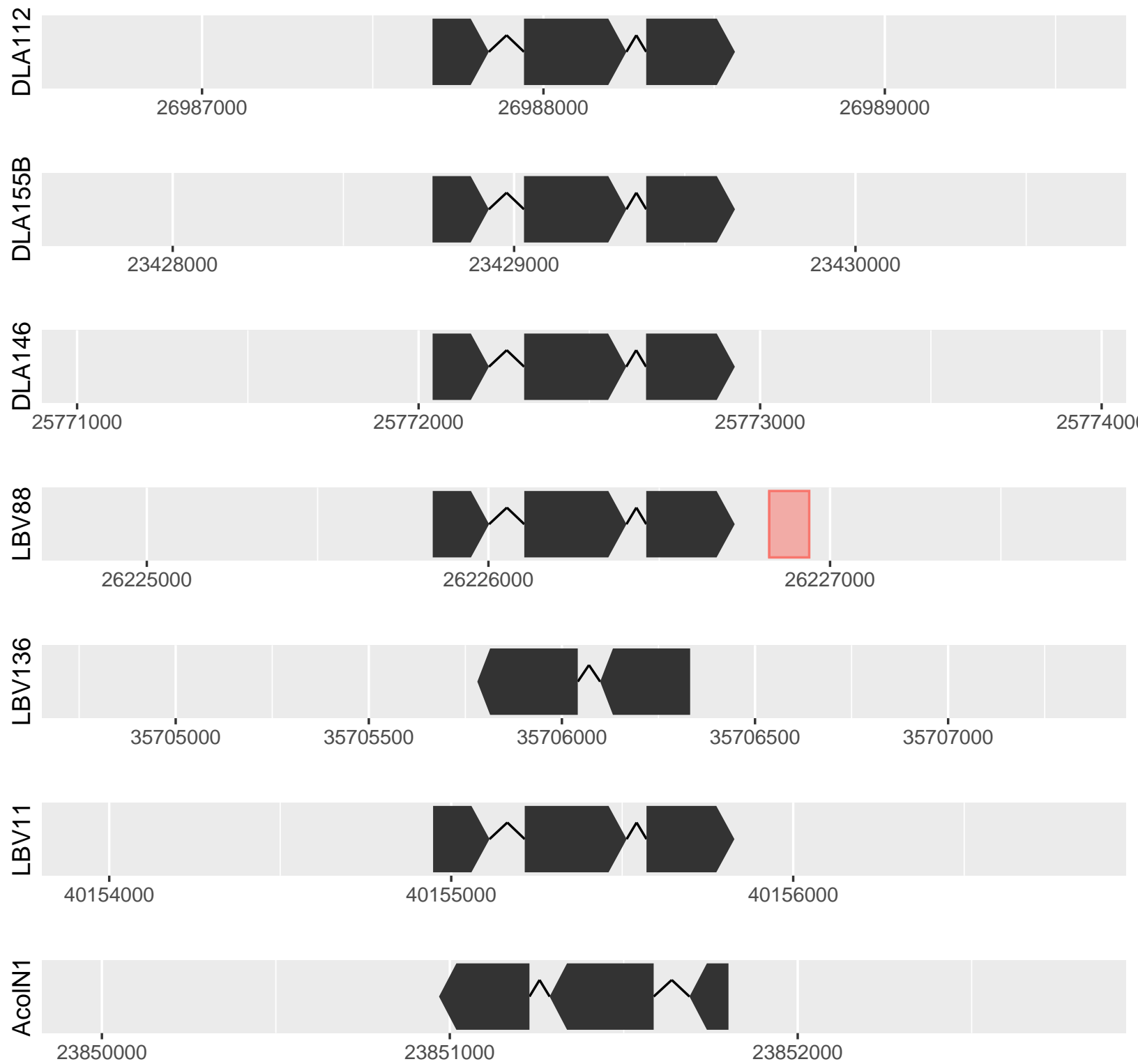
AGAP000741



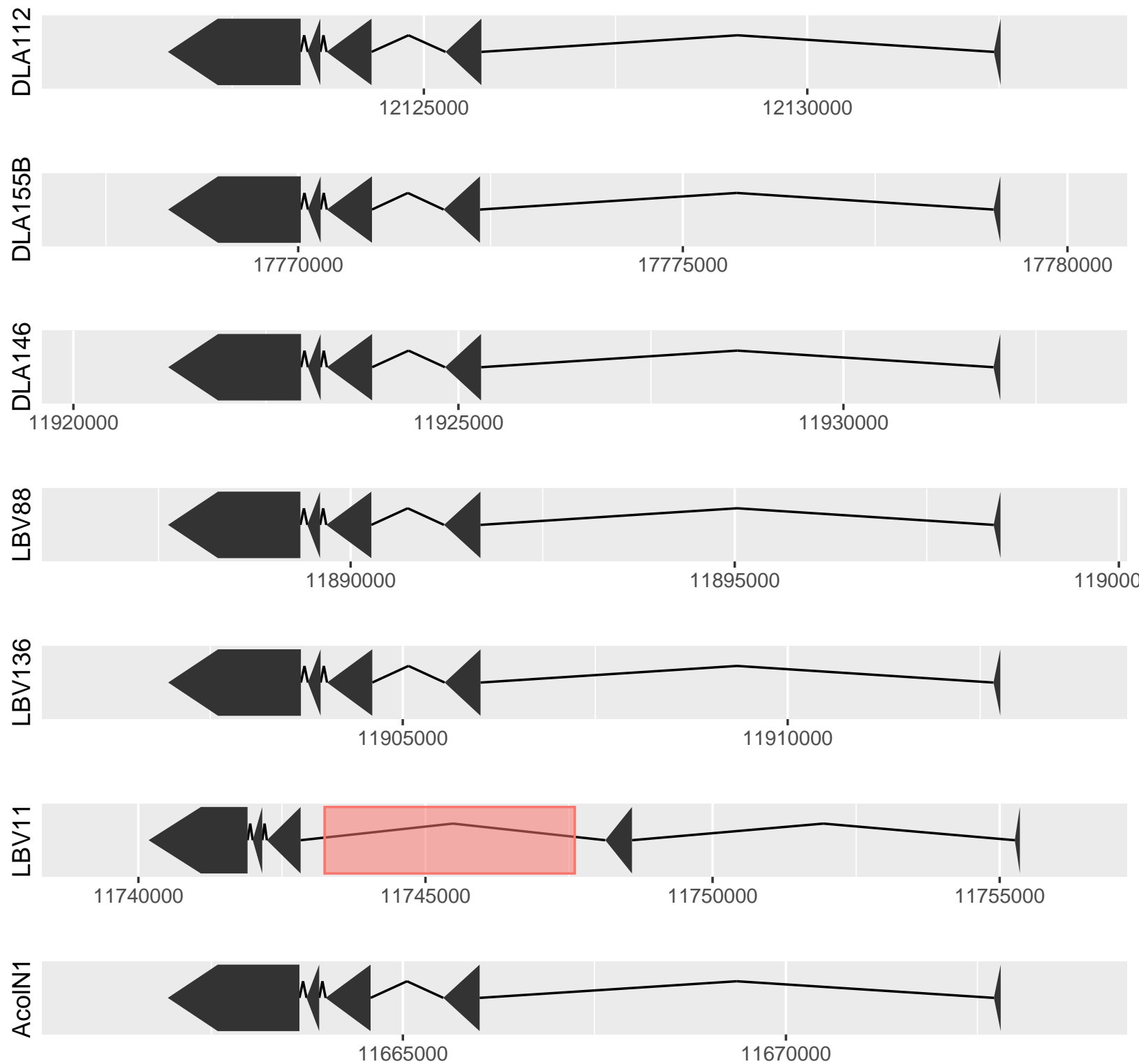
AGAP000894



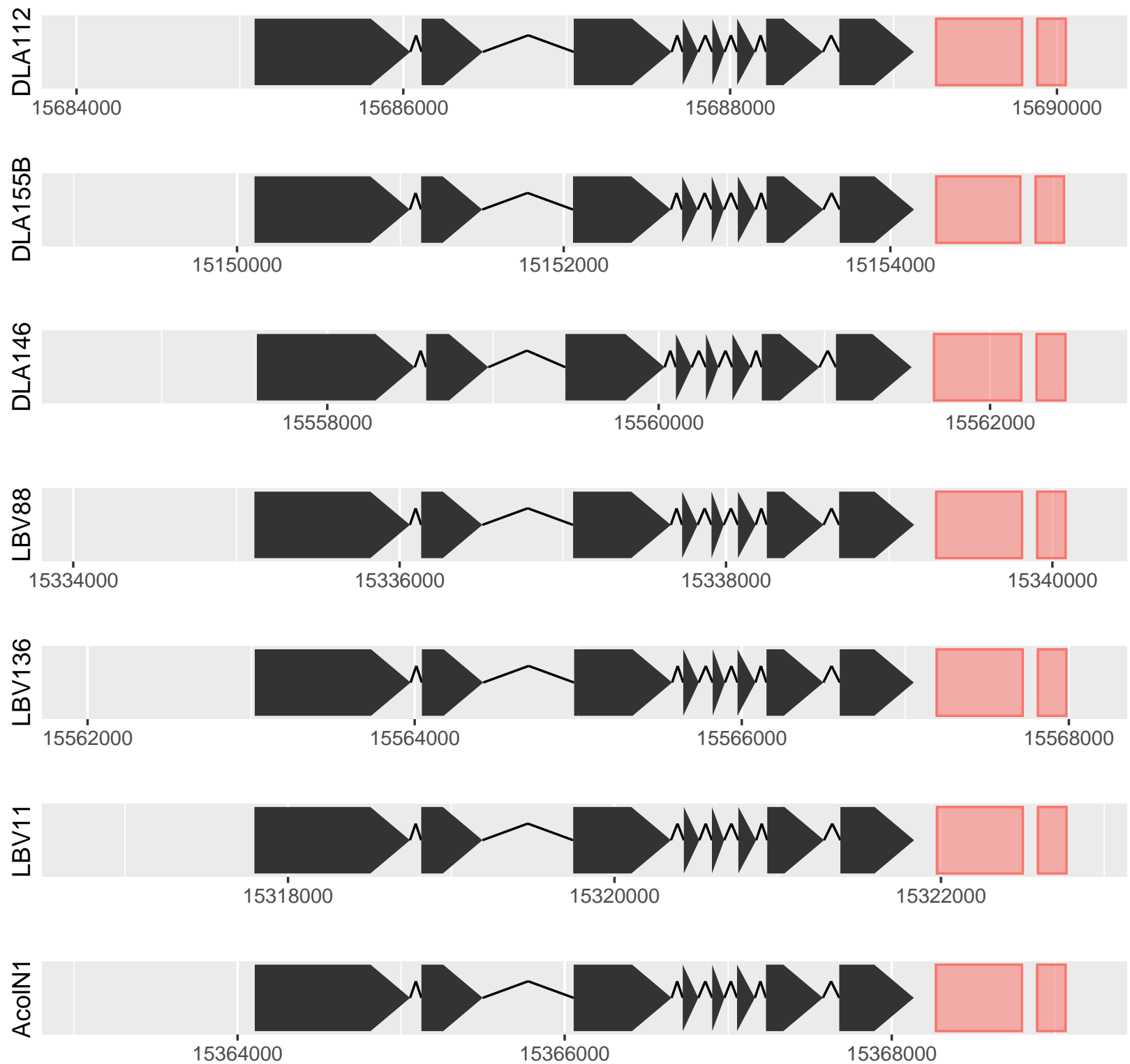
AGAP001079



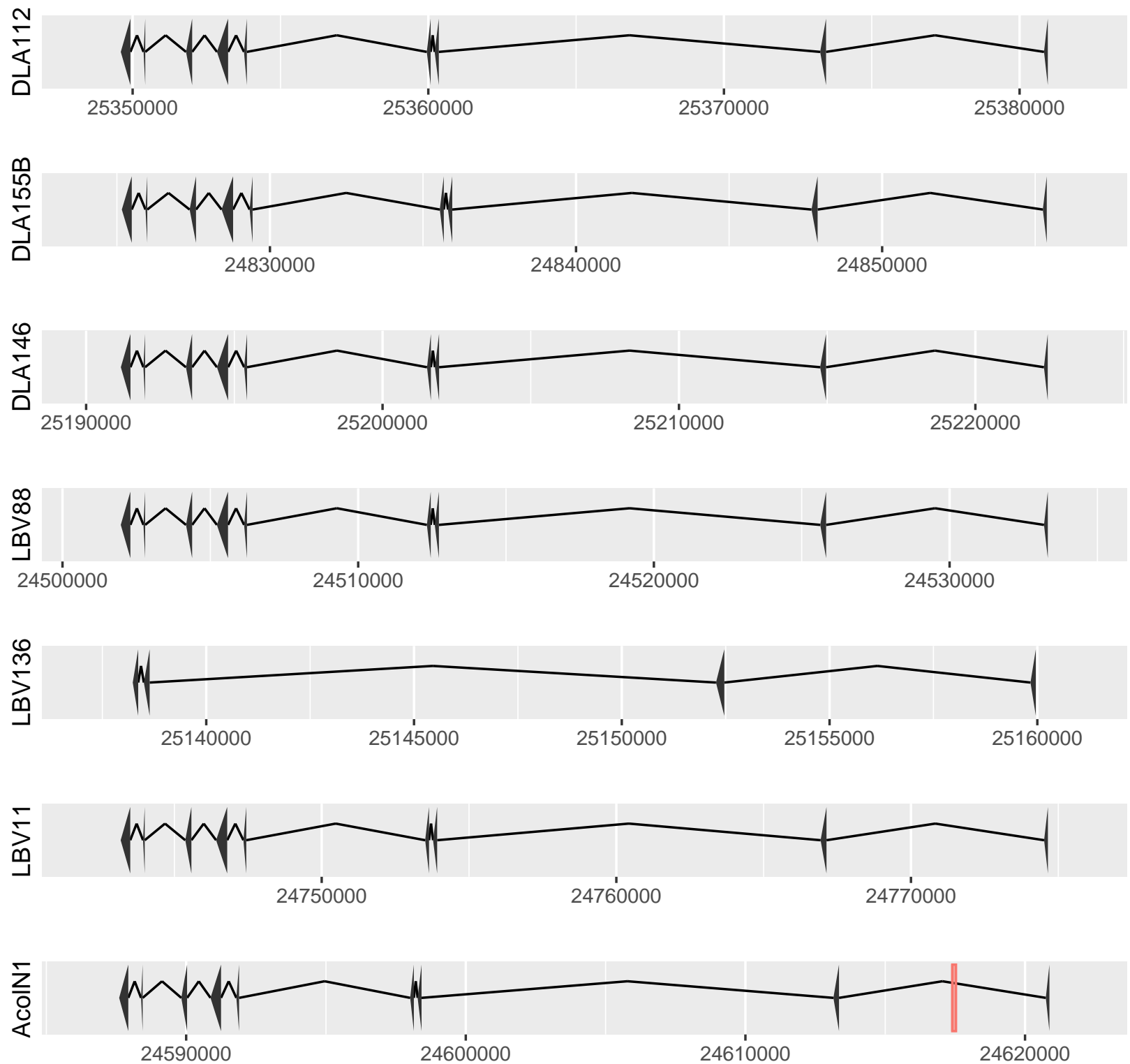
AGAP001882



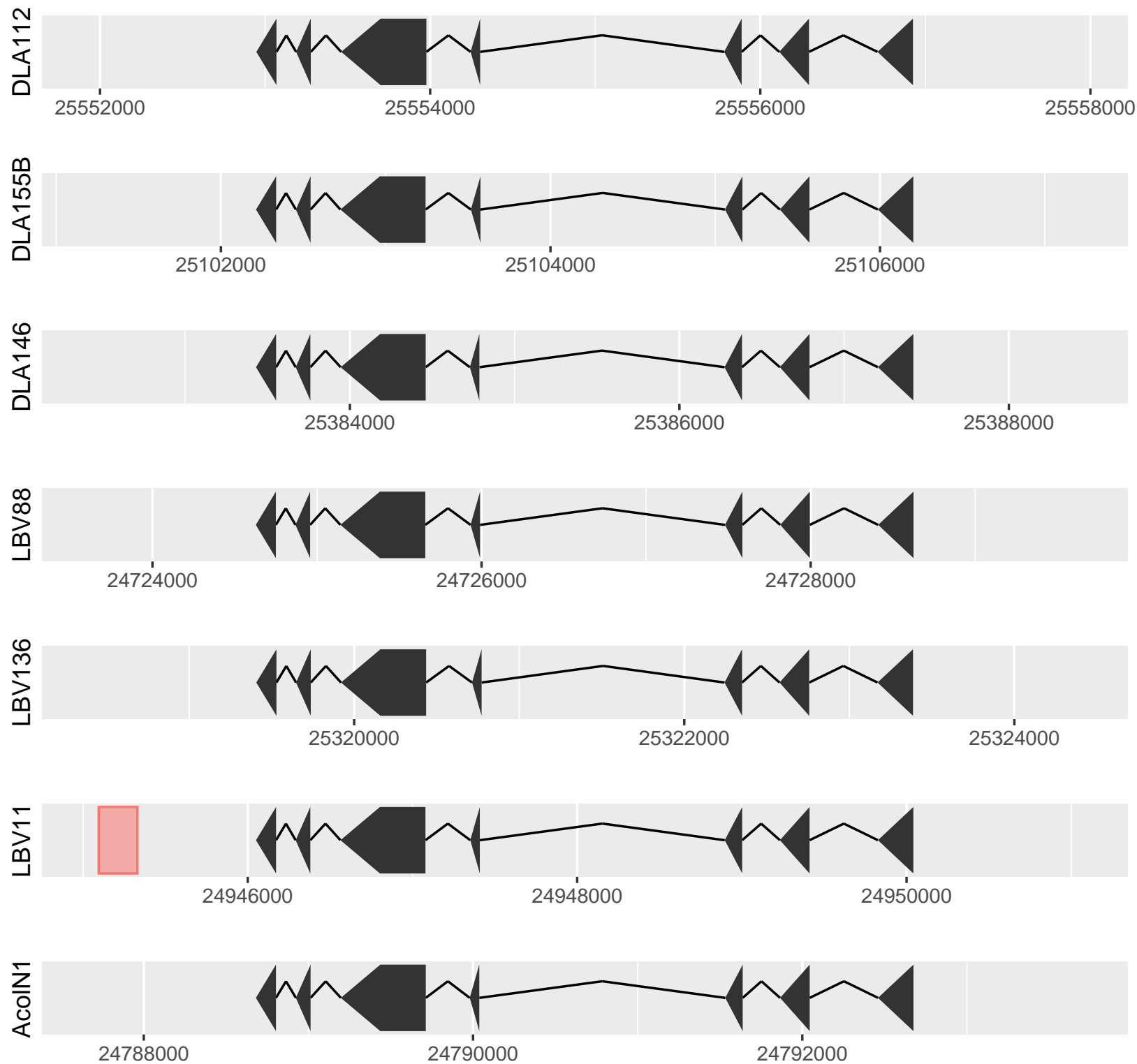
AGAP002100



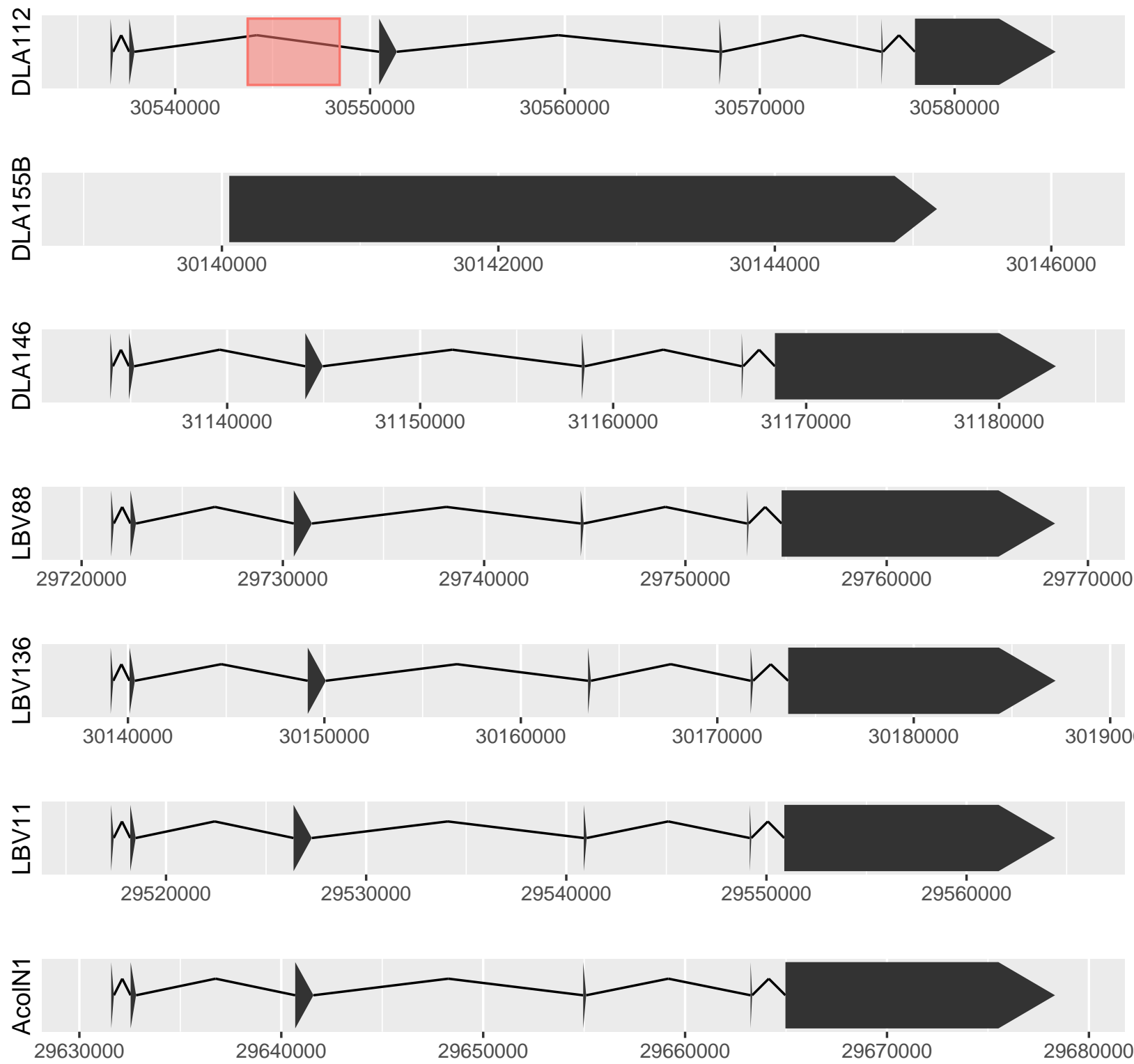
AGAP002628



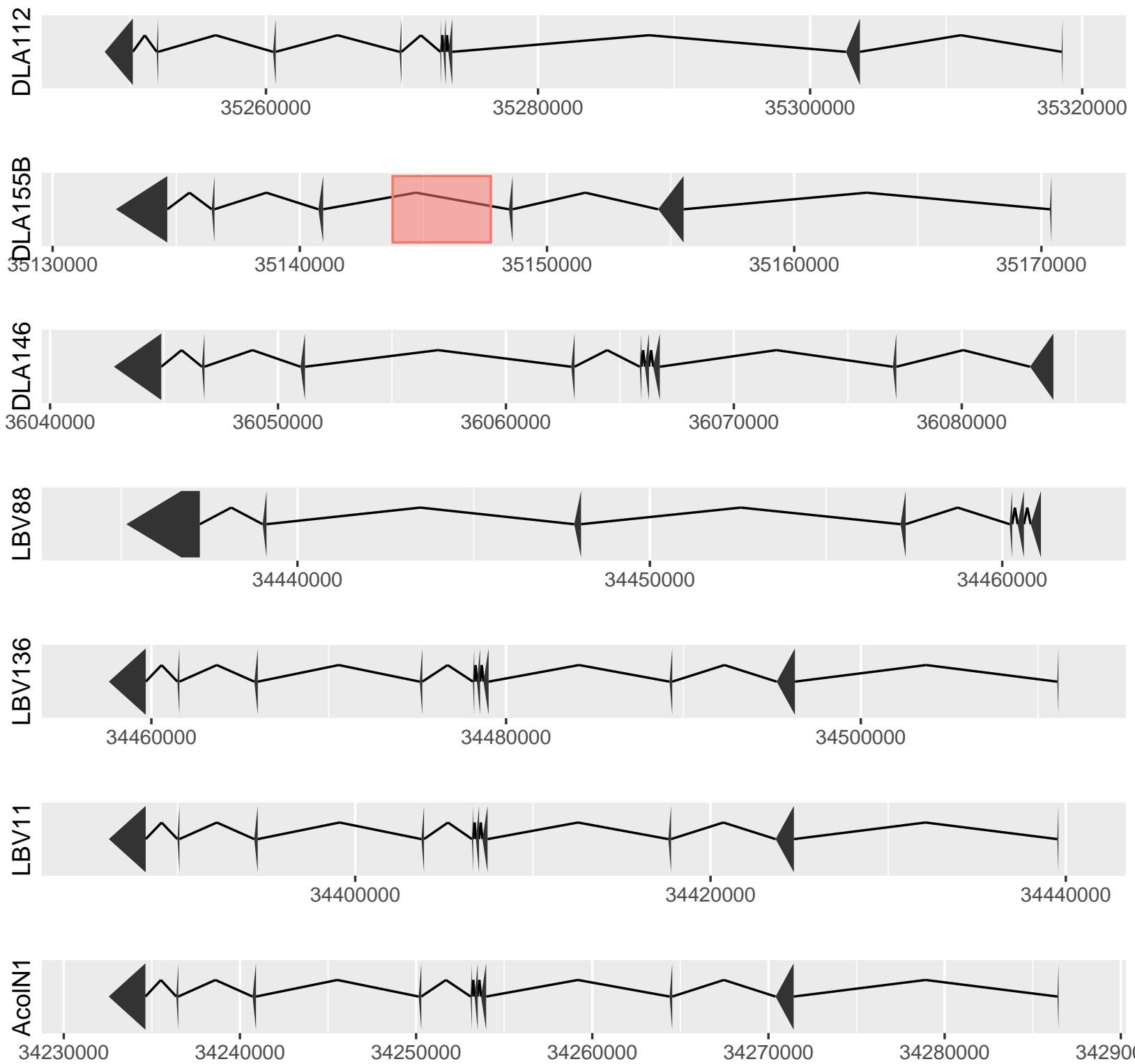
AGAP002633



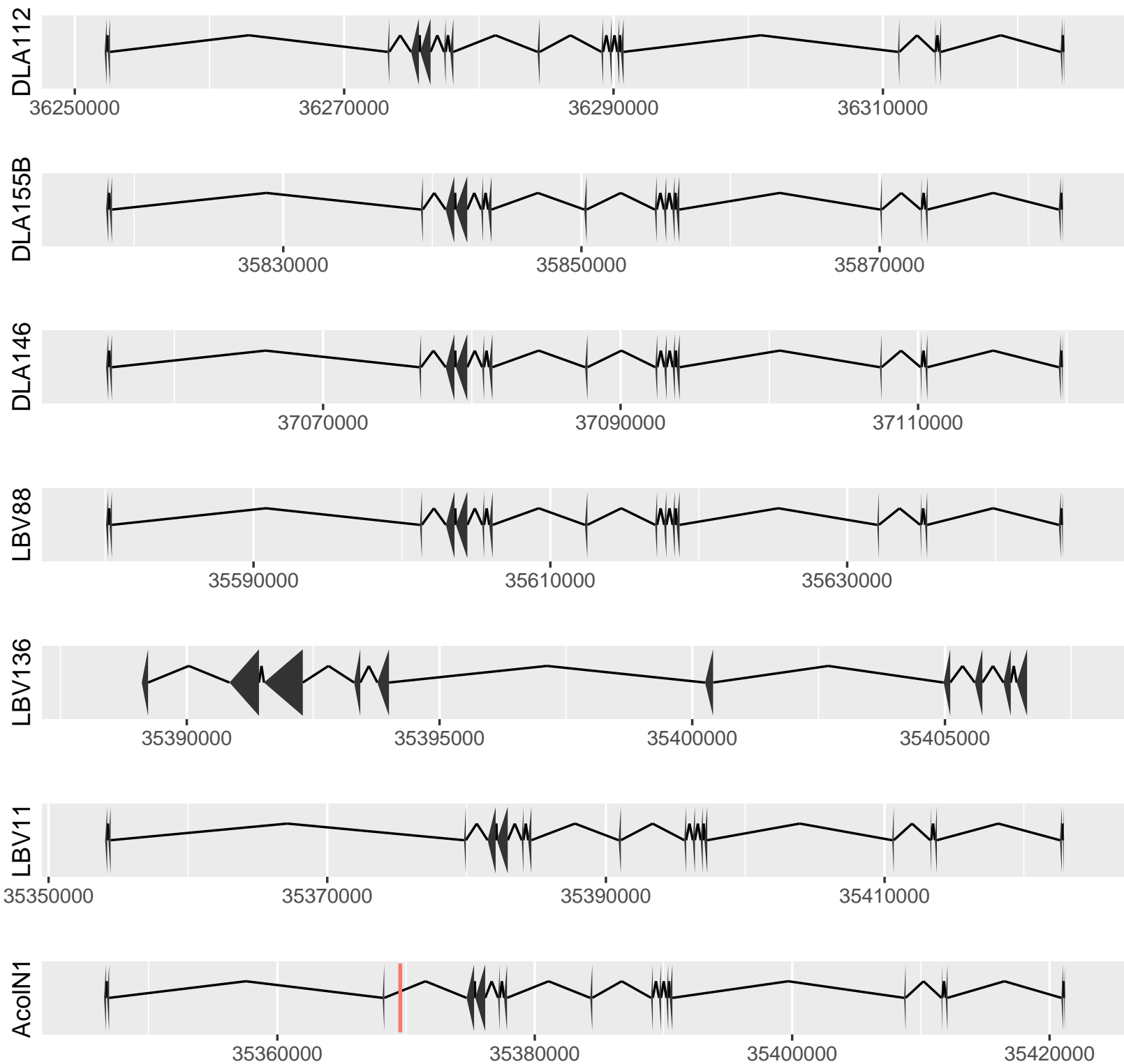
AGAP002916



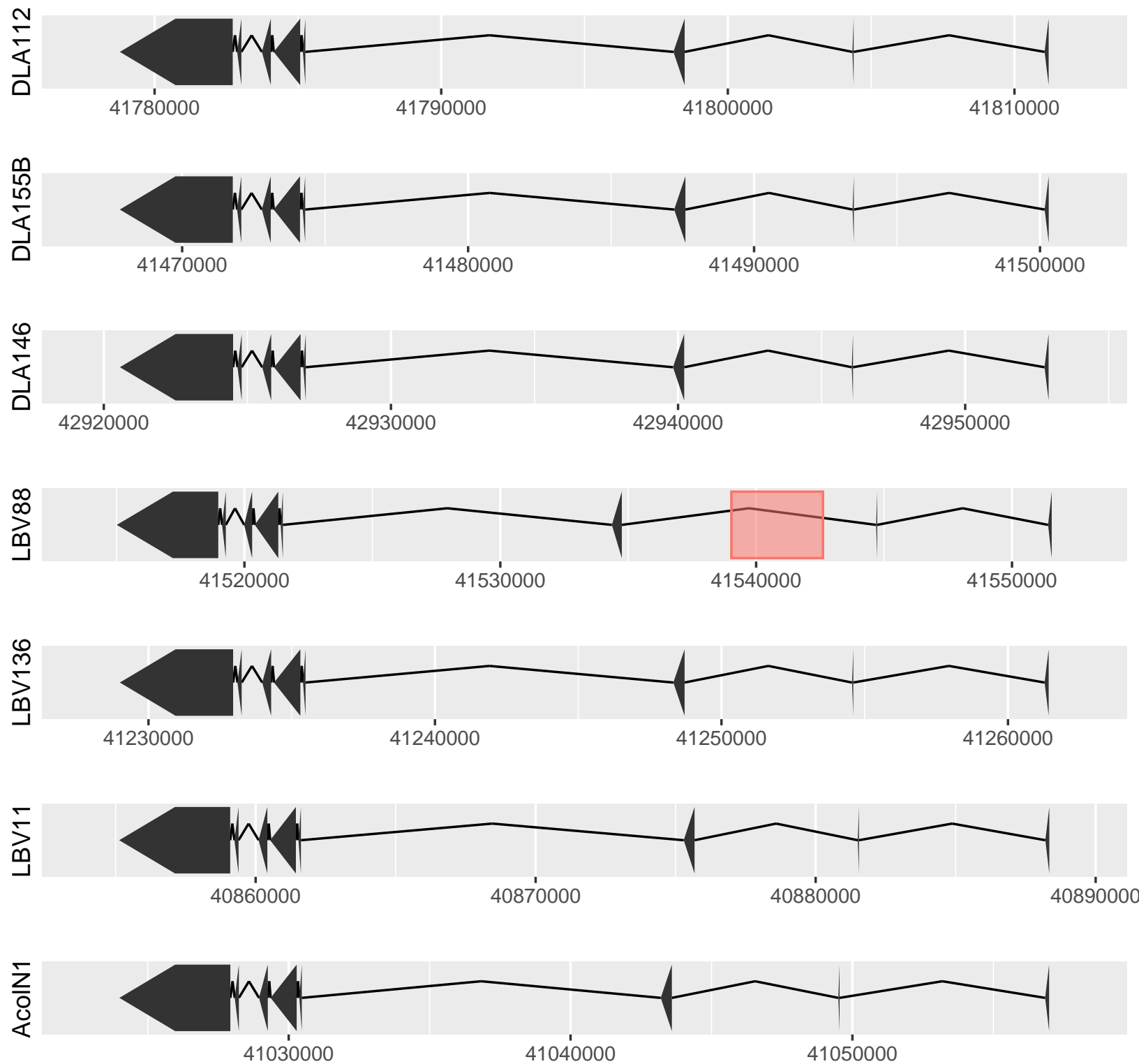
AGAP003244



AGAP003305



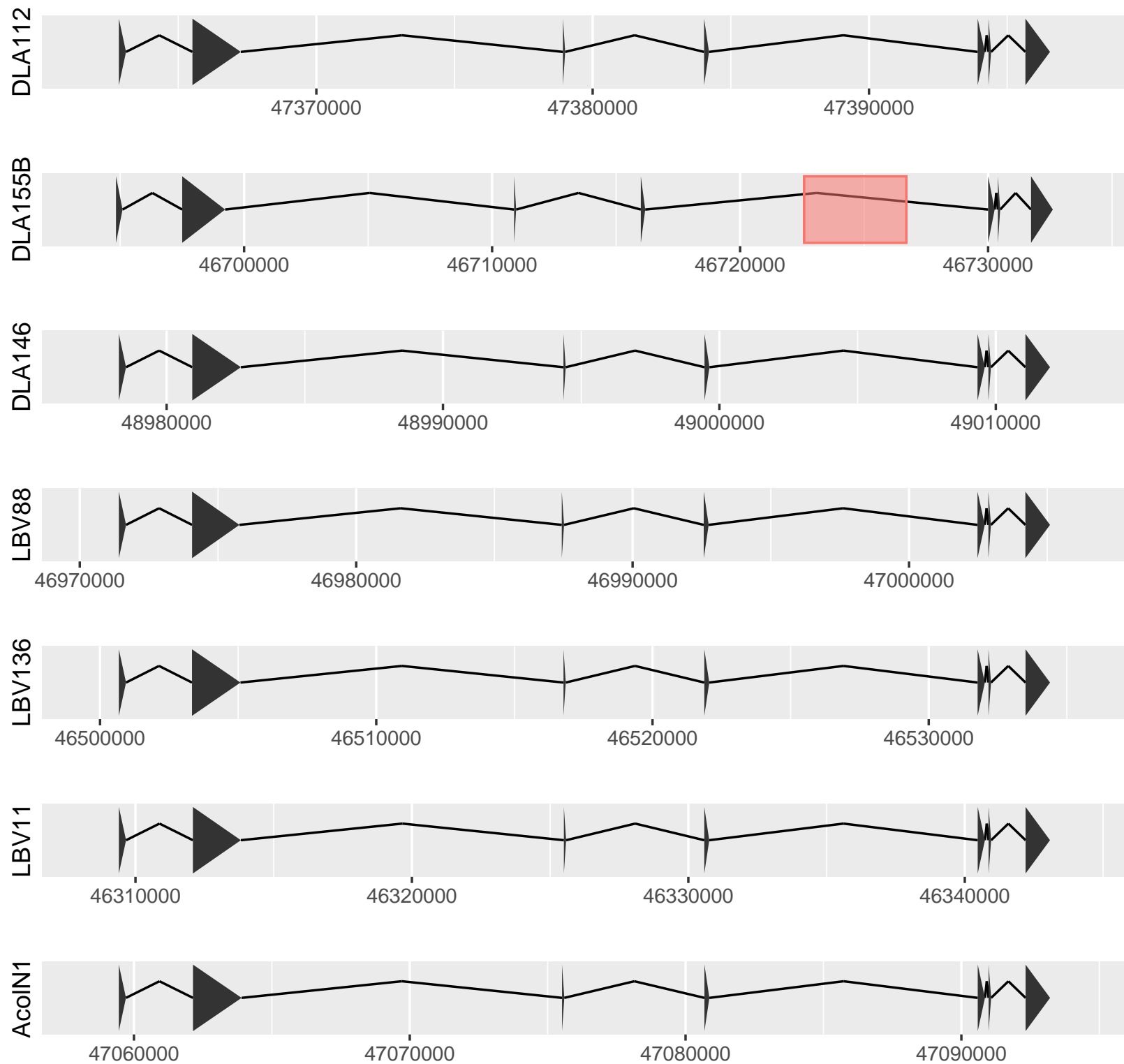
AGAP003649



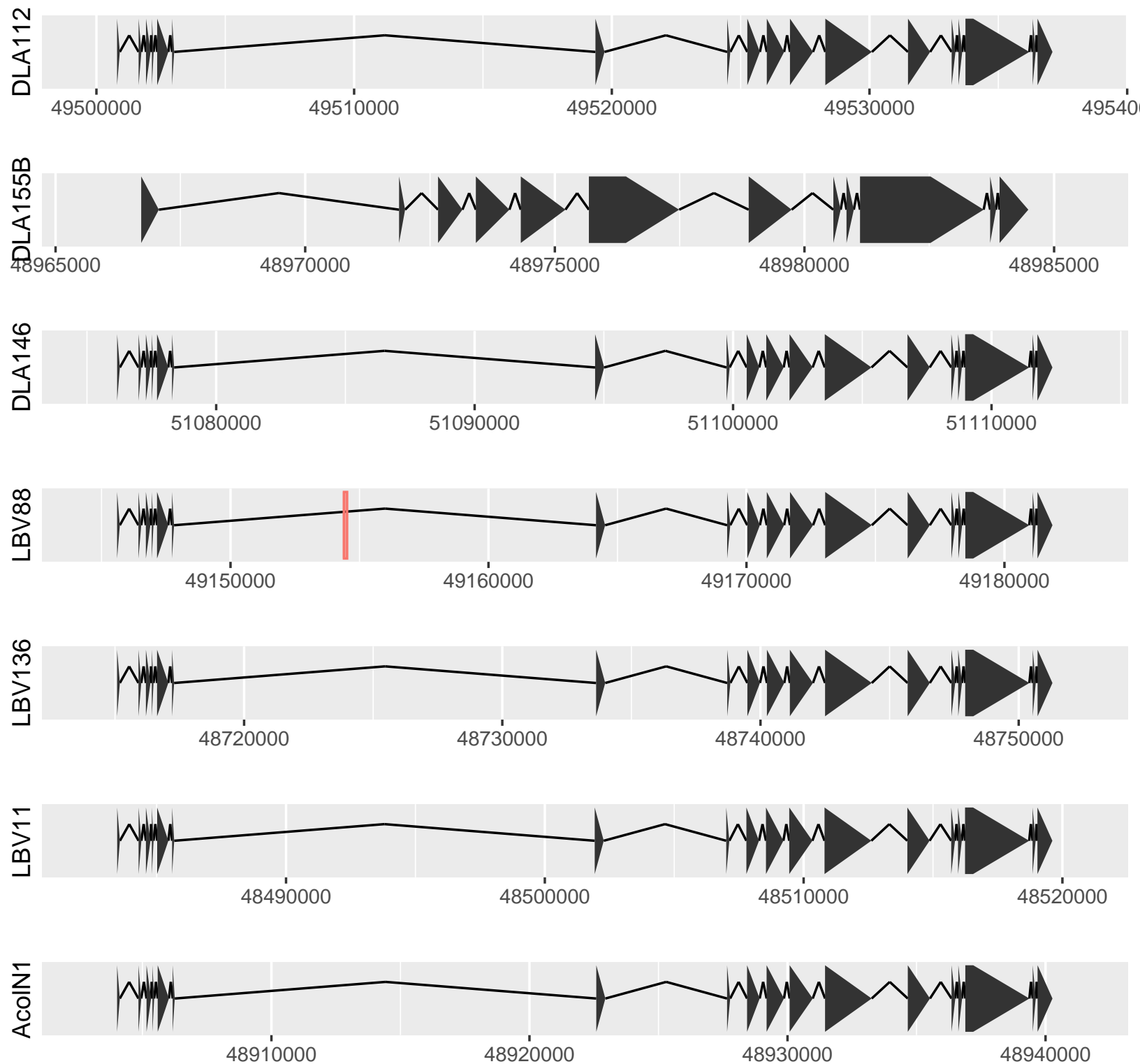
AGAP003676



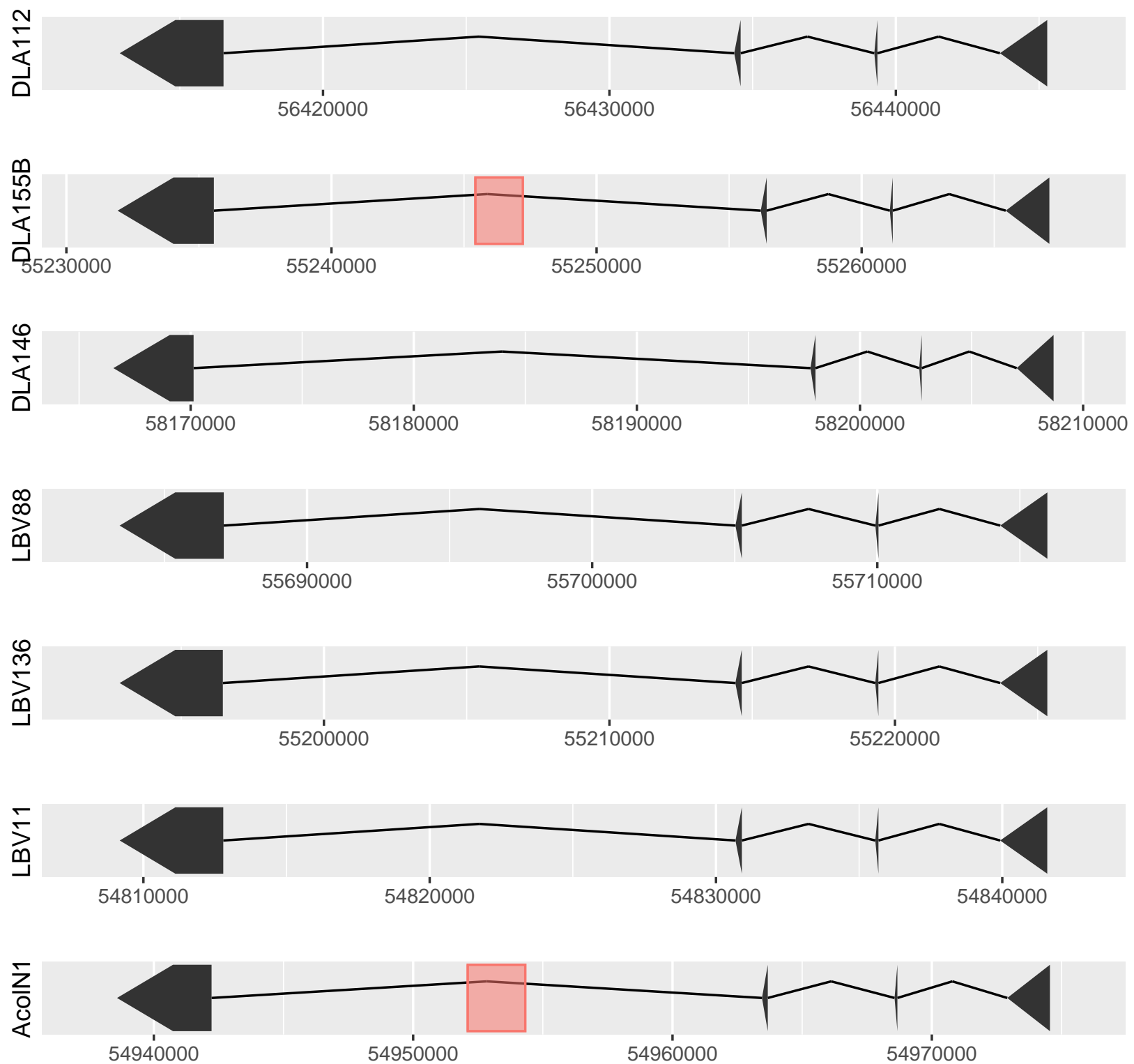
AGAP003931



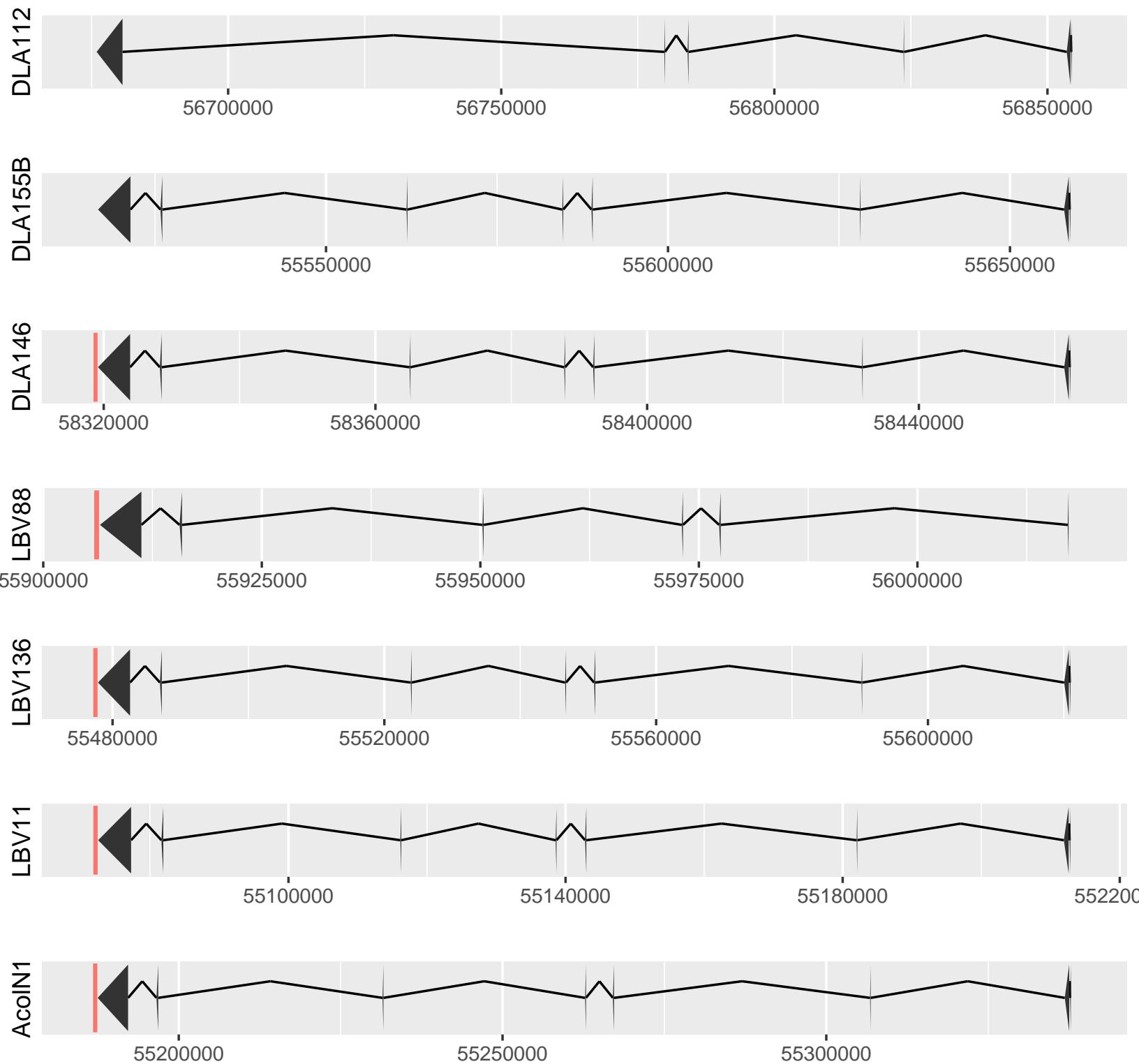
AGAP004039



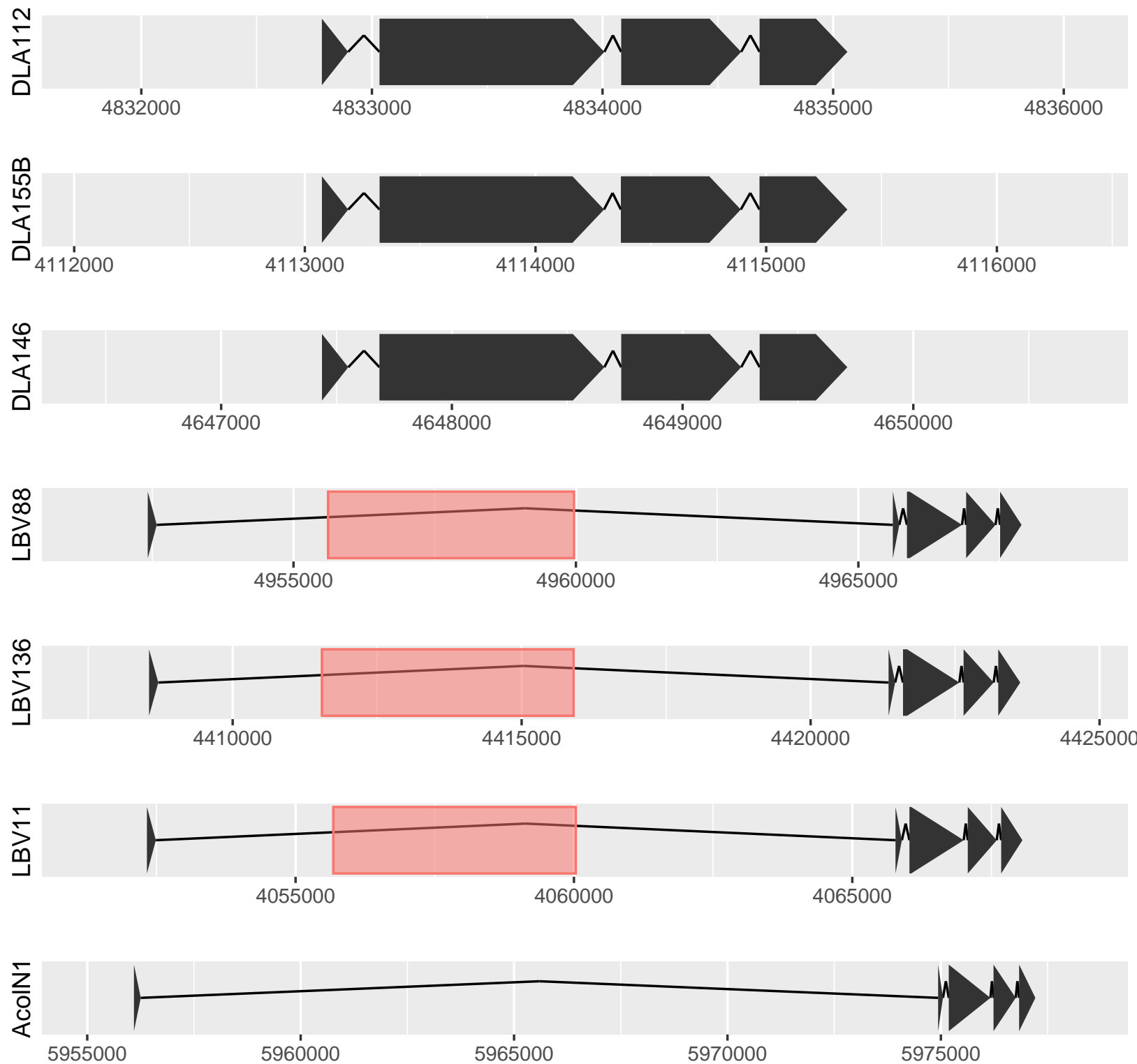
AGAP004353



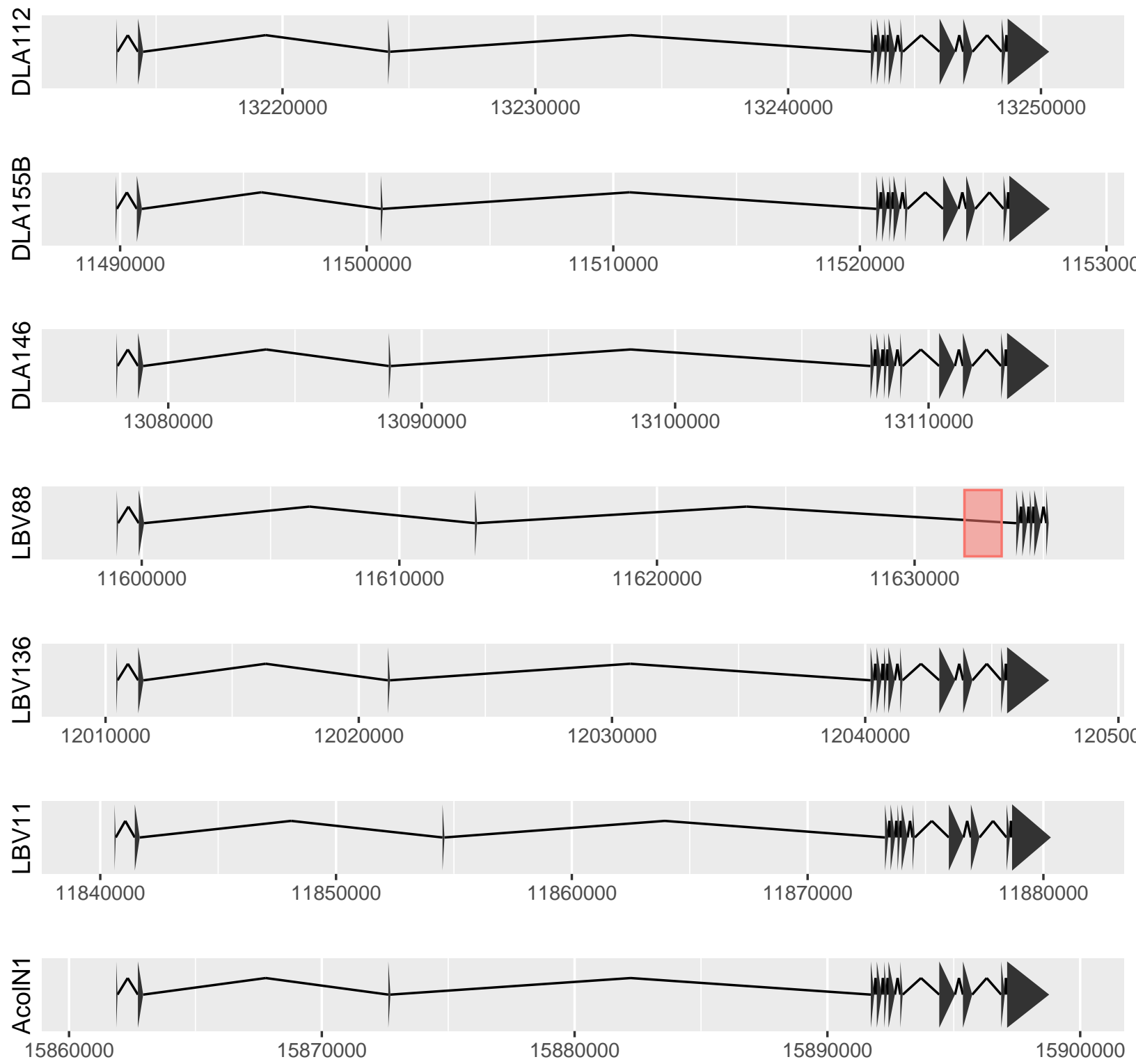
AGAP004369



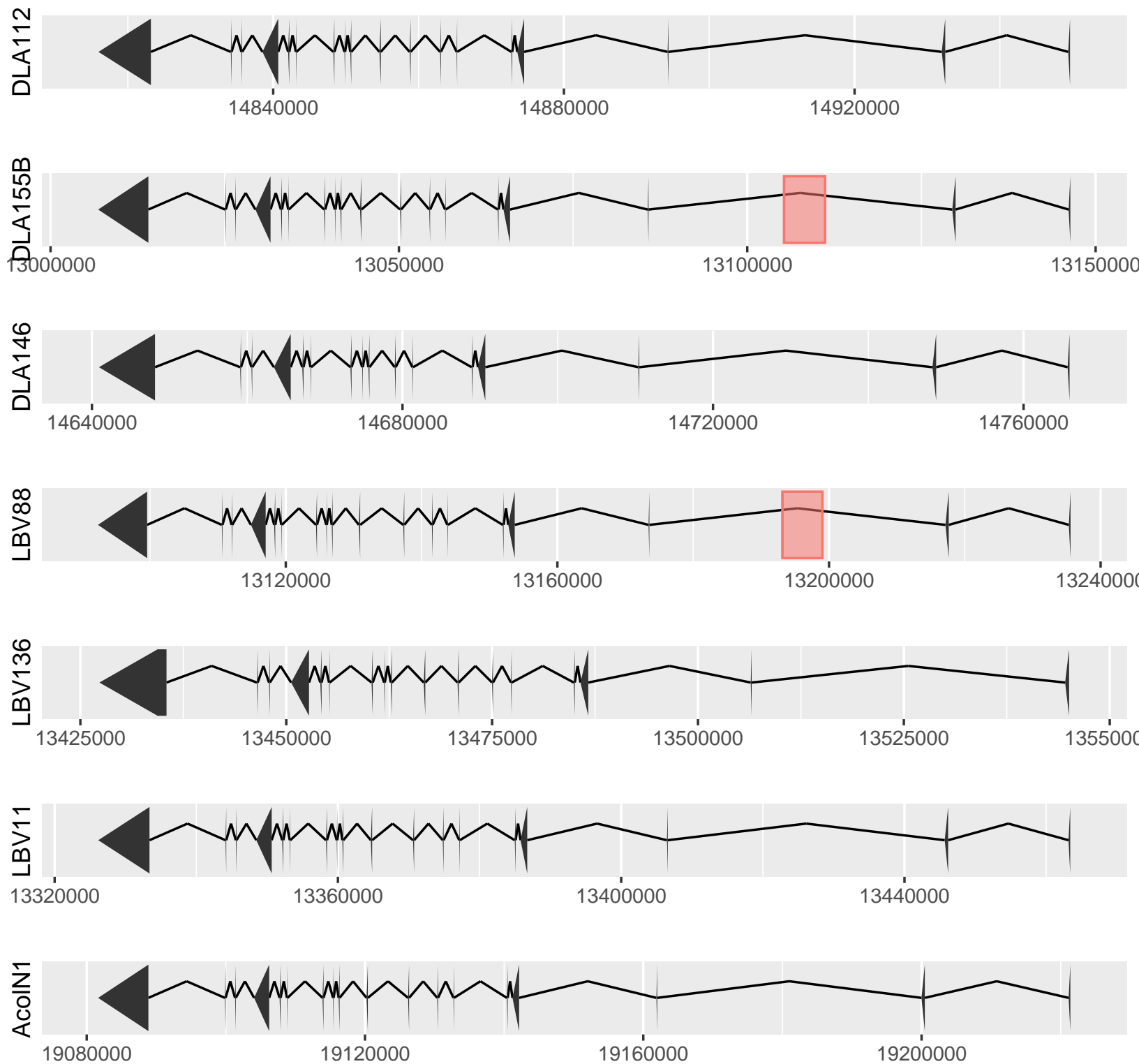
AGAP004724



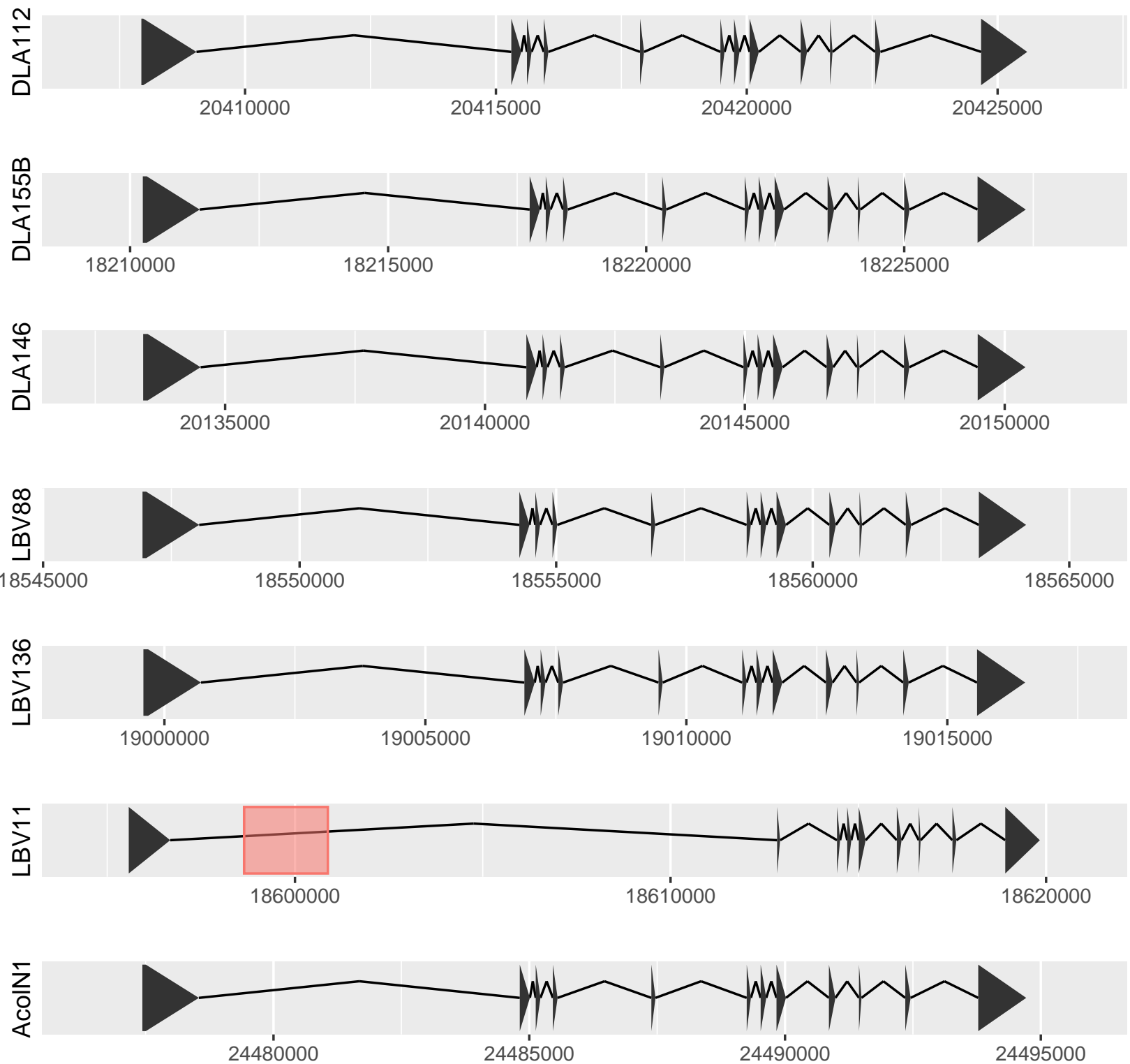
AGAP005076



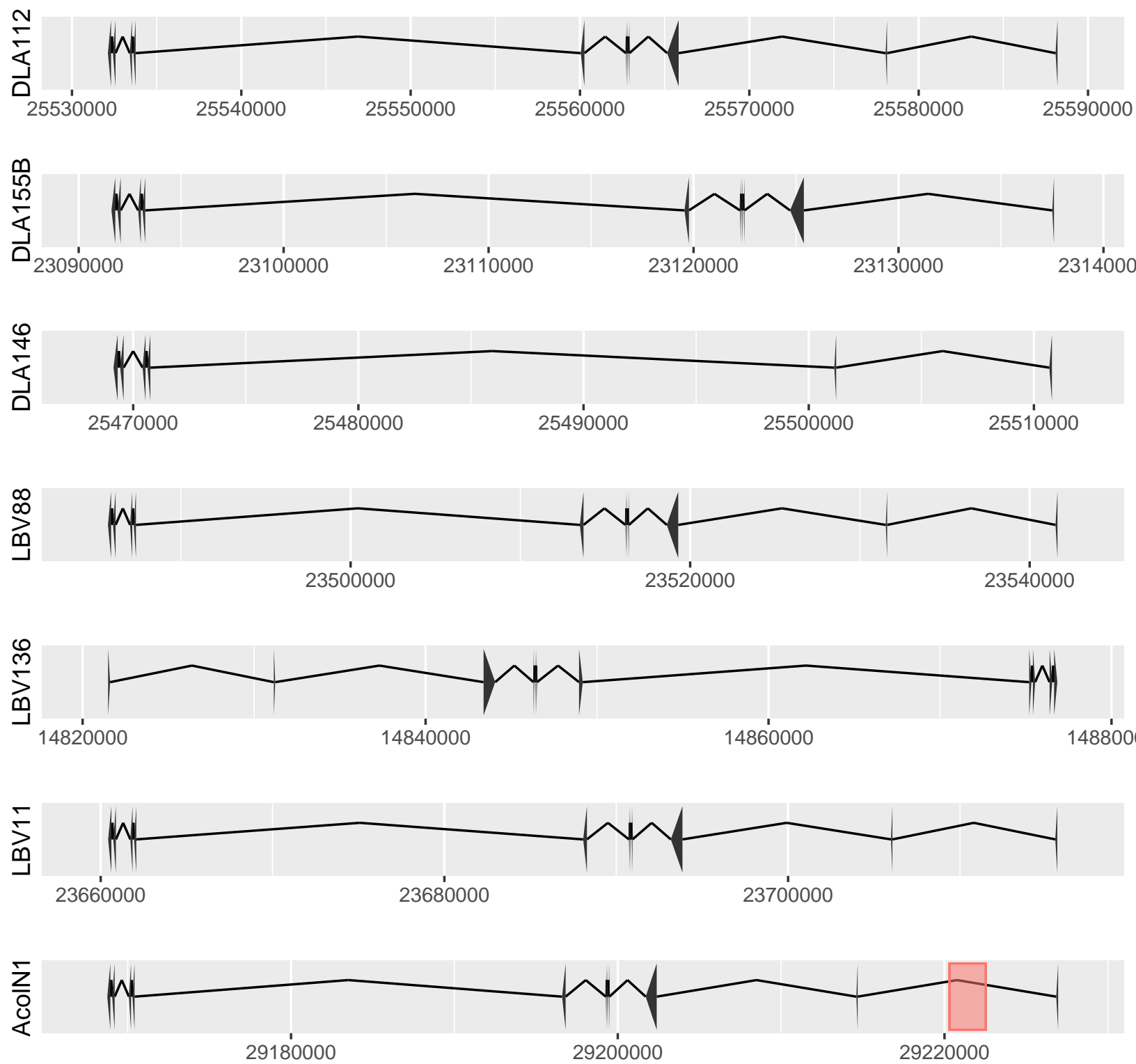
AGAP005165



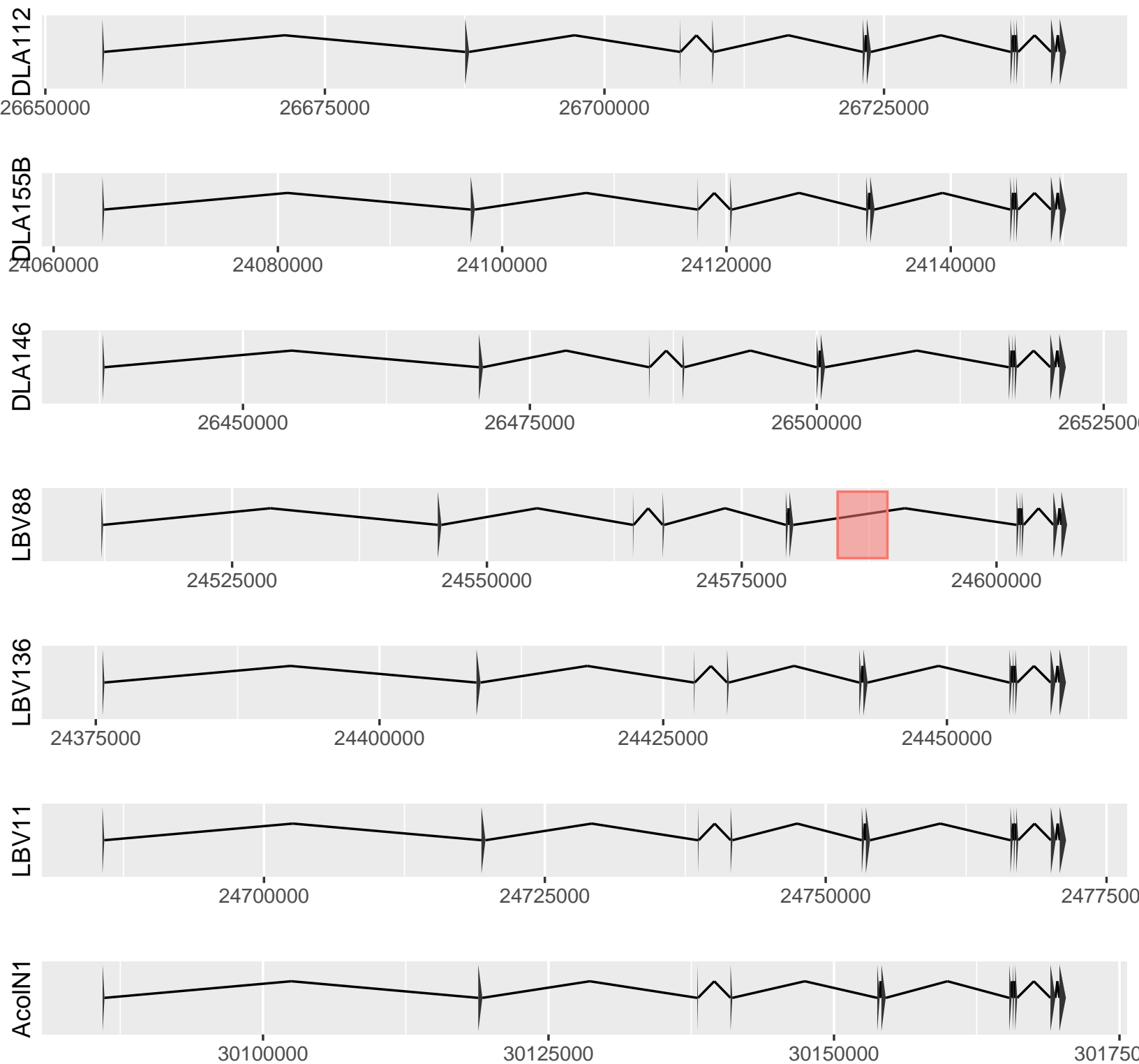
AGAP005464



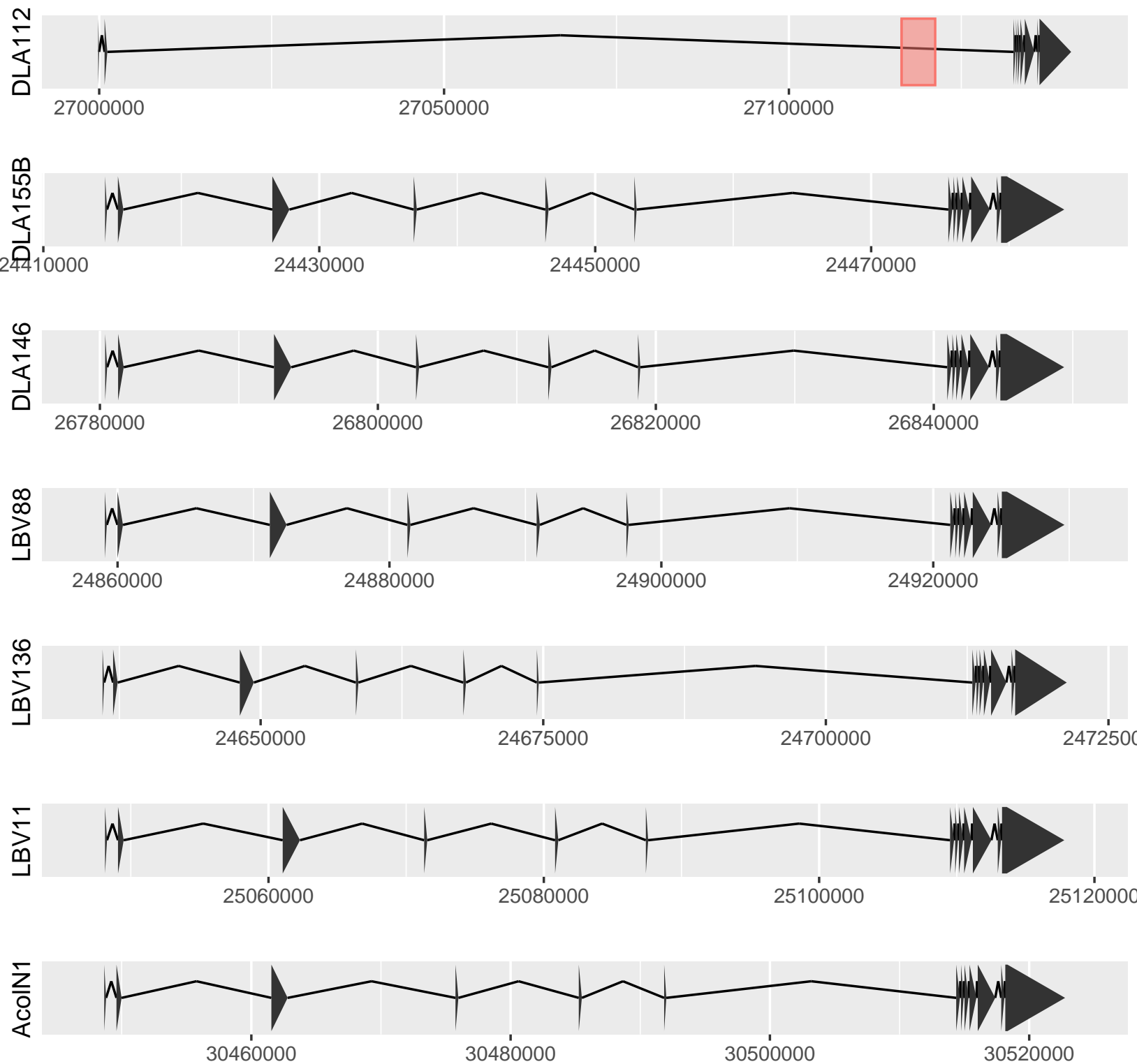
AGAP005773



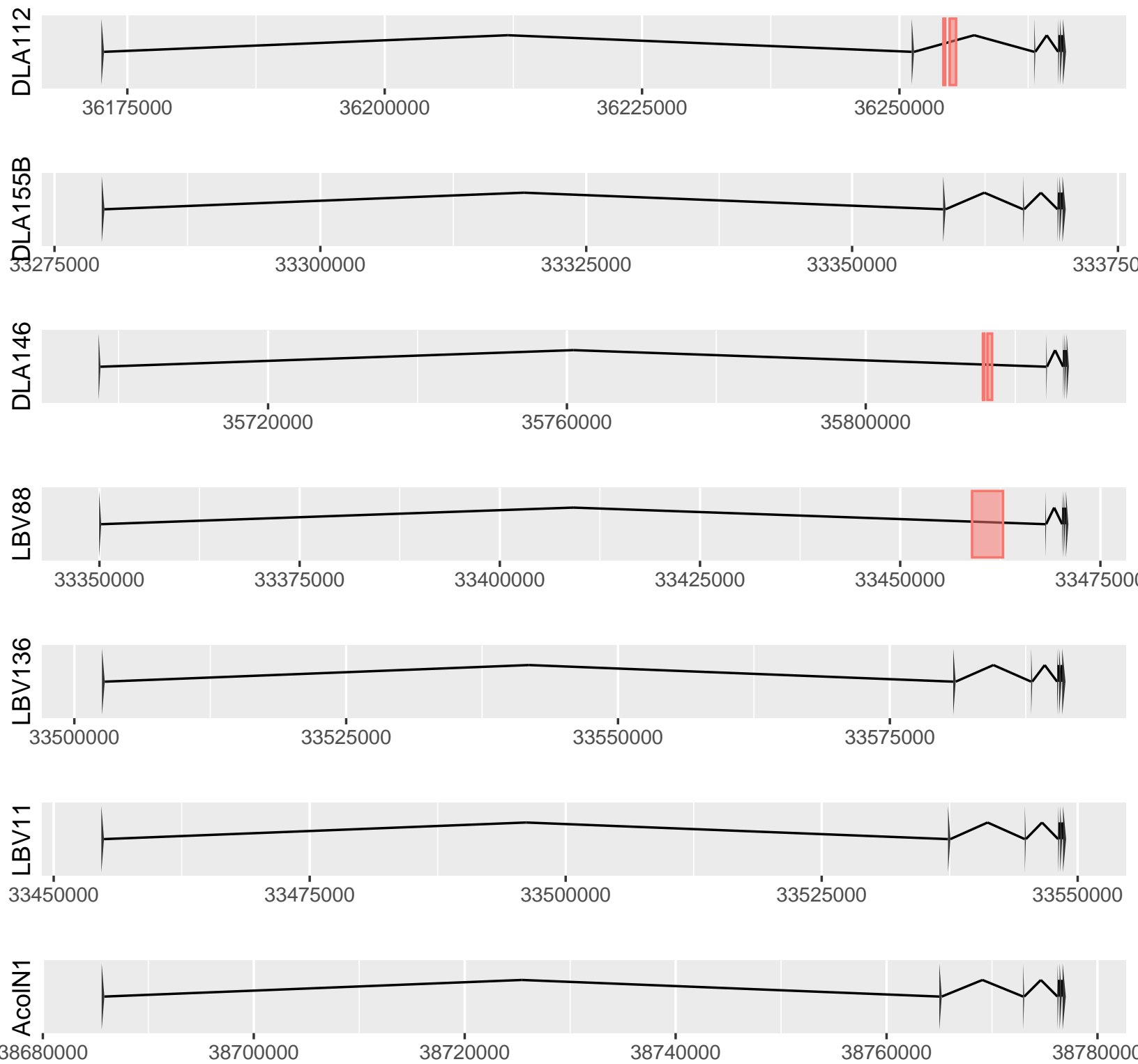
AGAP005806



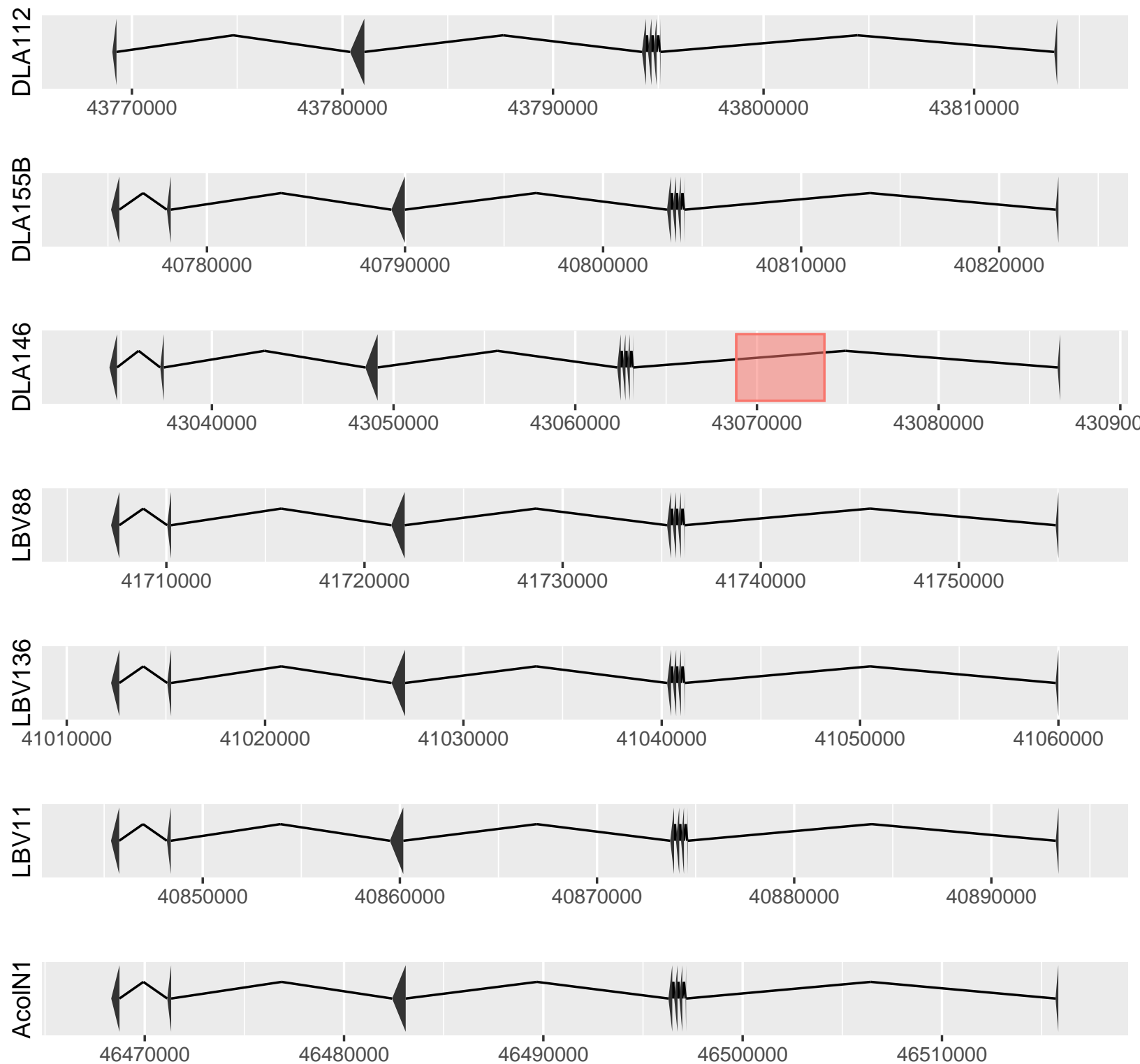
AGAP005816



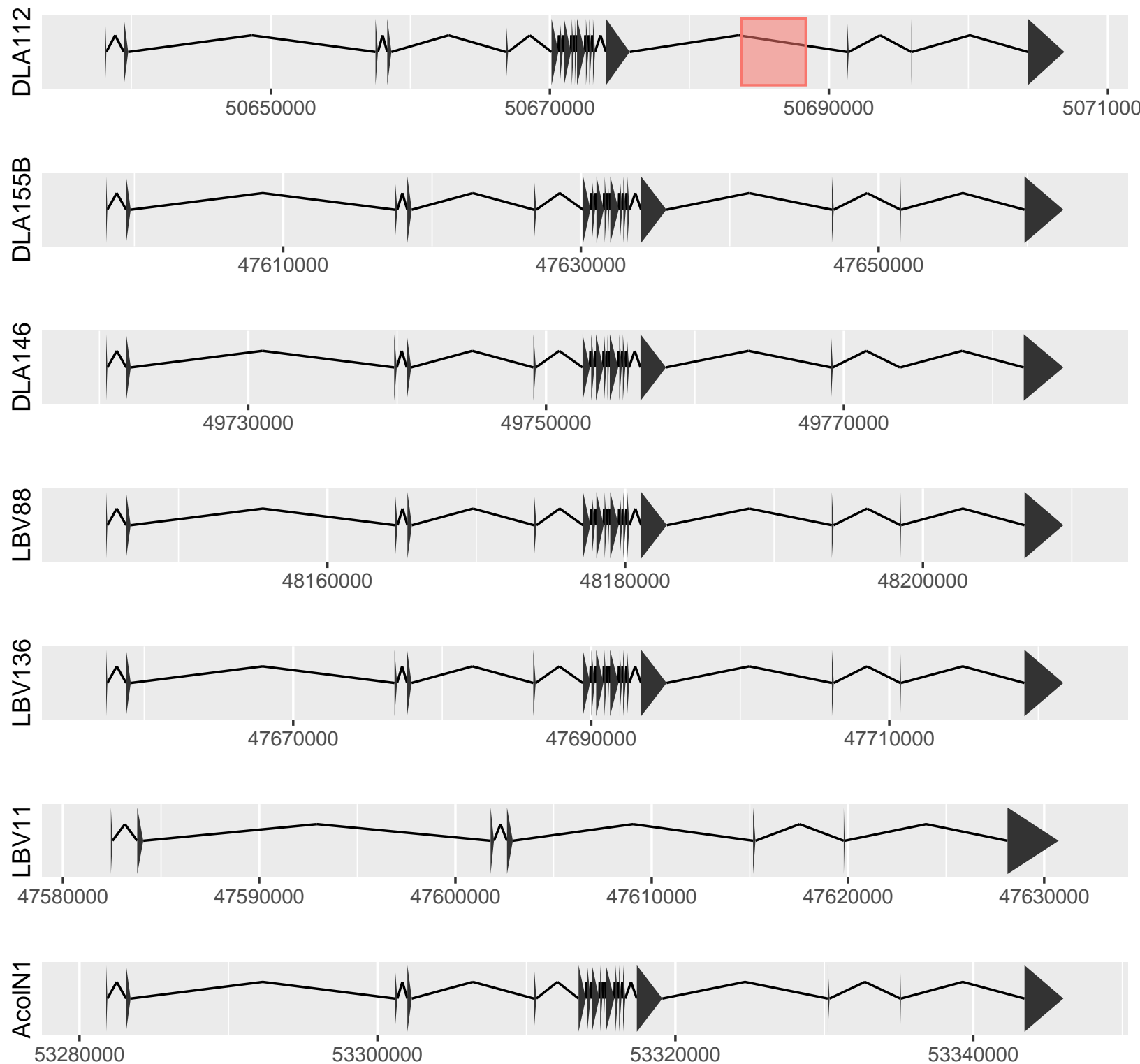
AGAP006330



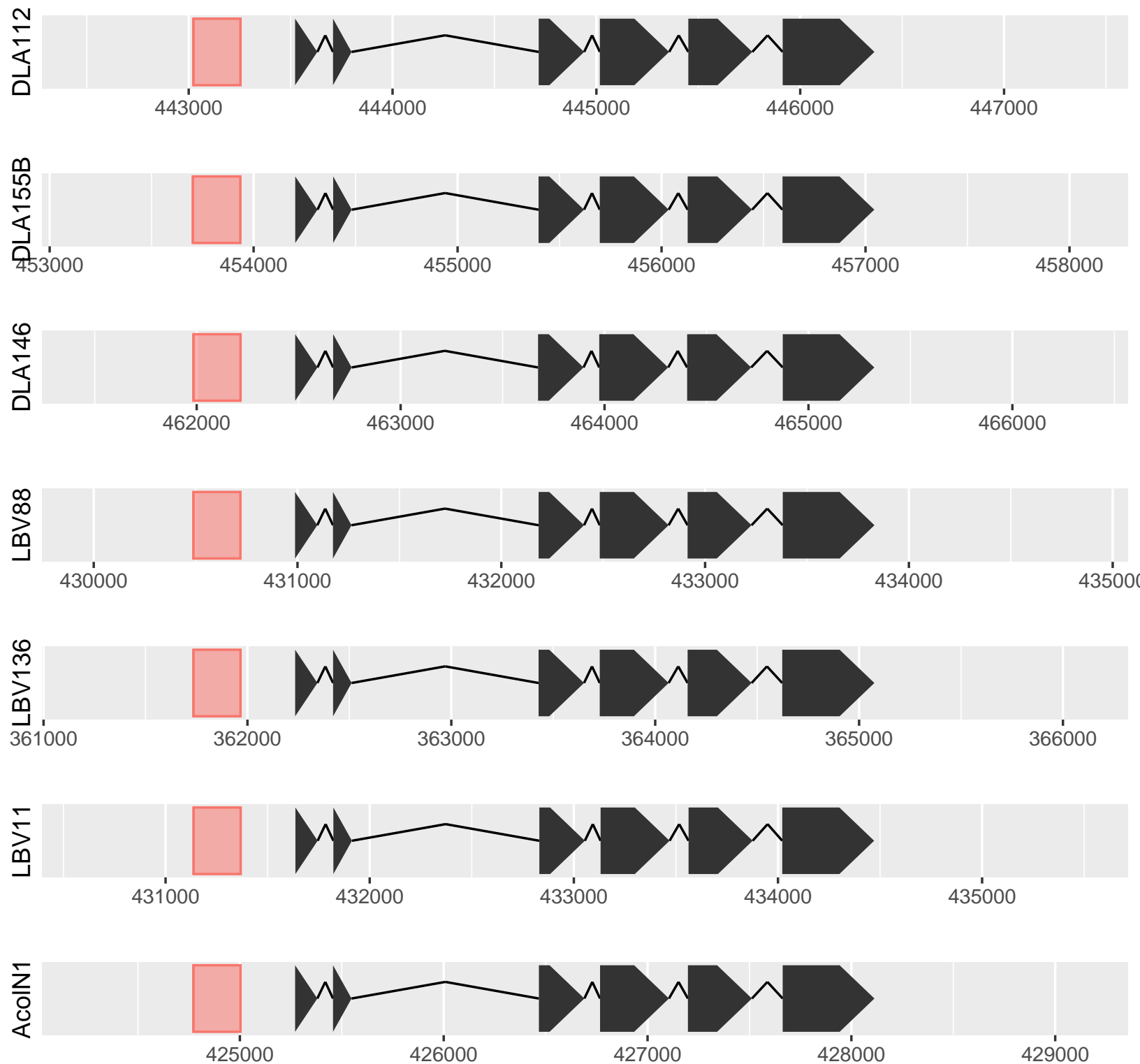
AGAP006721



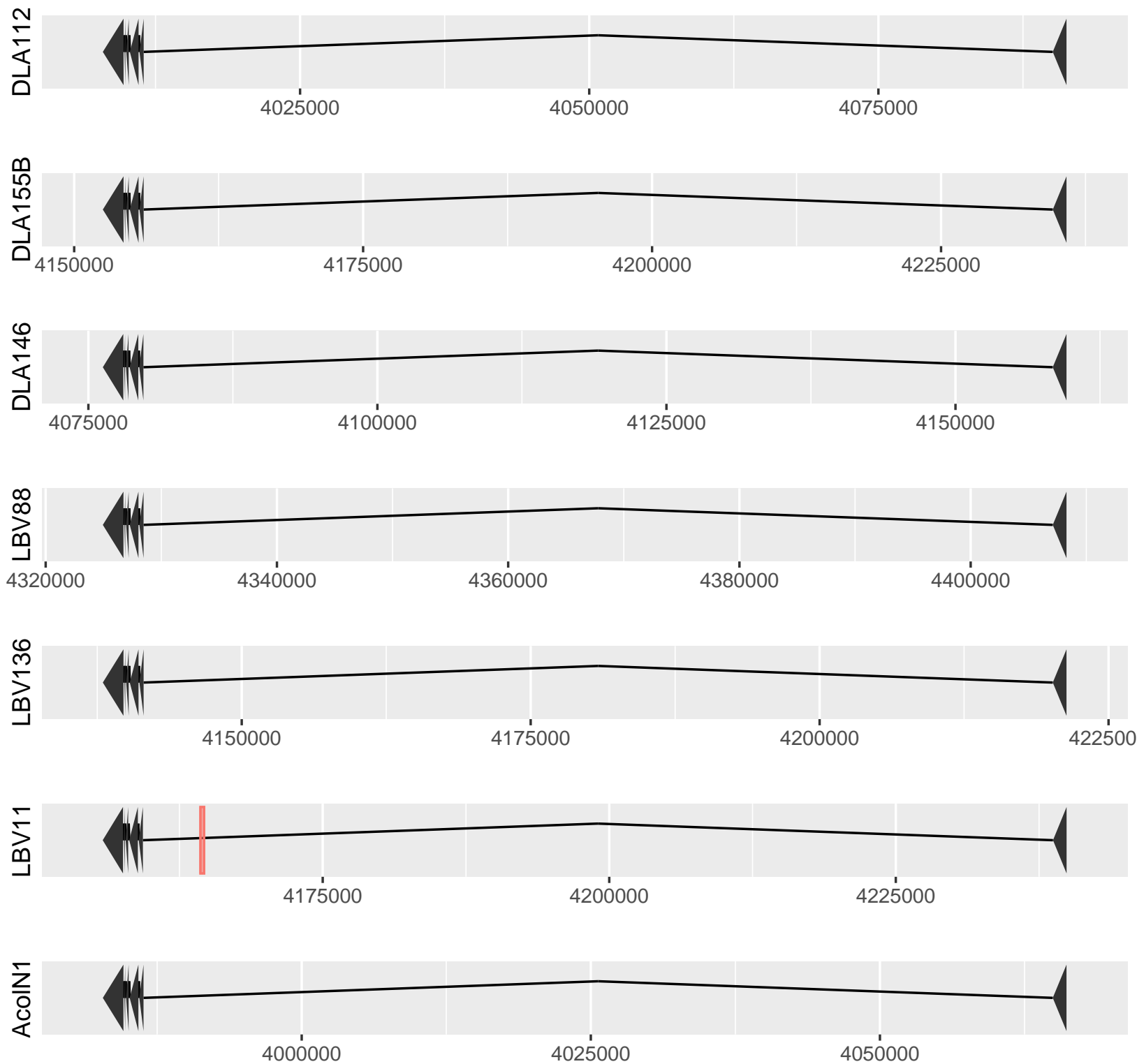
AGAP007137



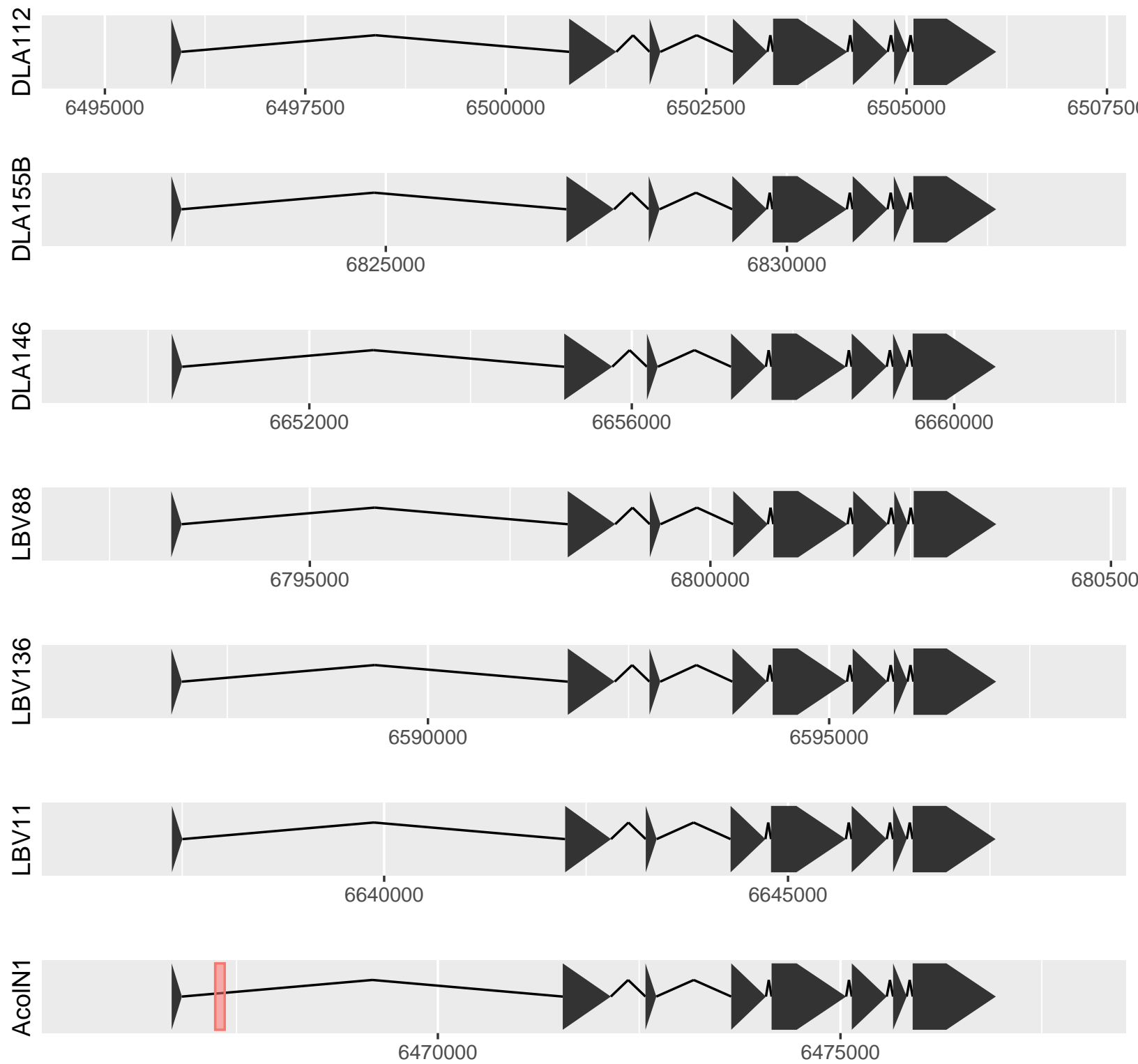
AGAP007760



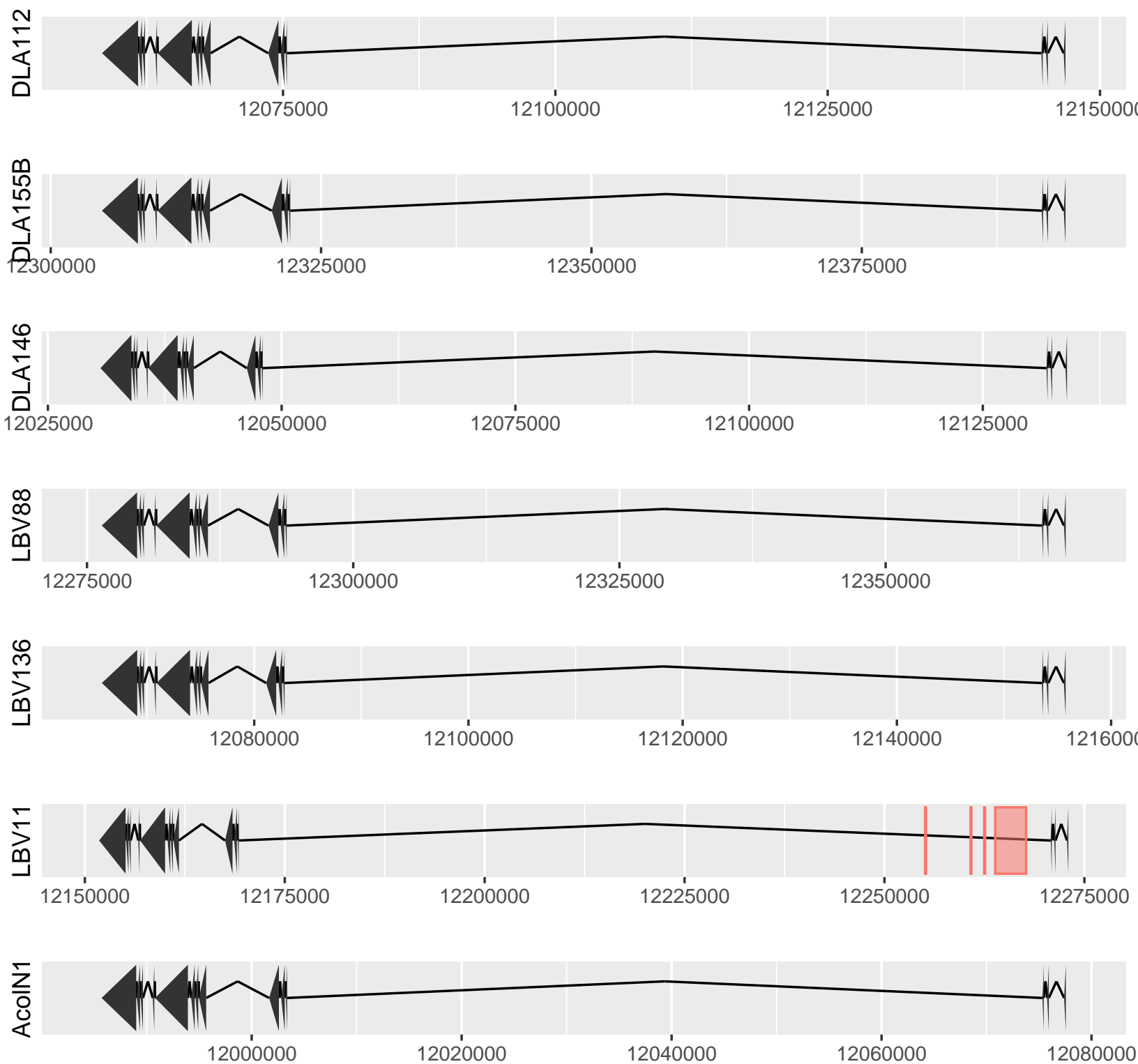
AGAP008009



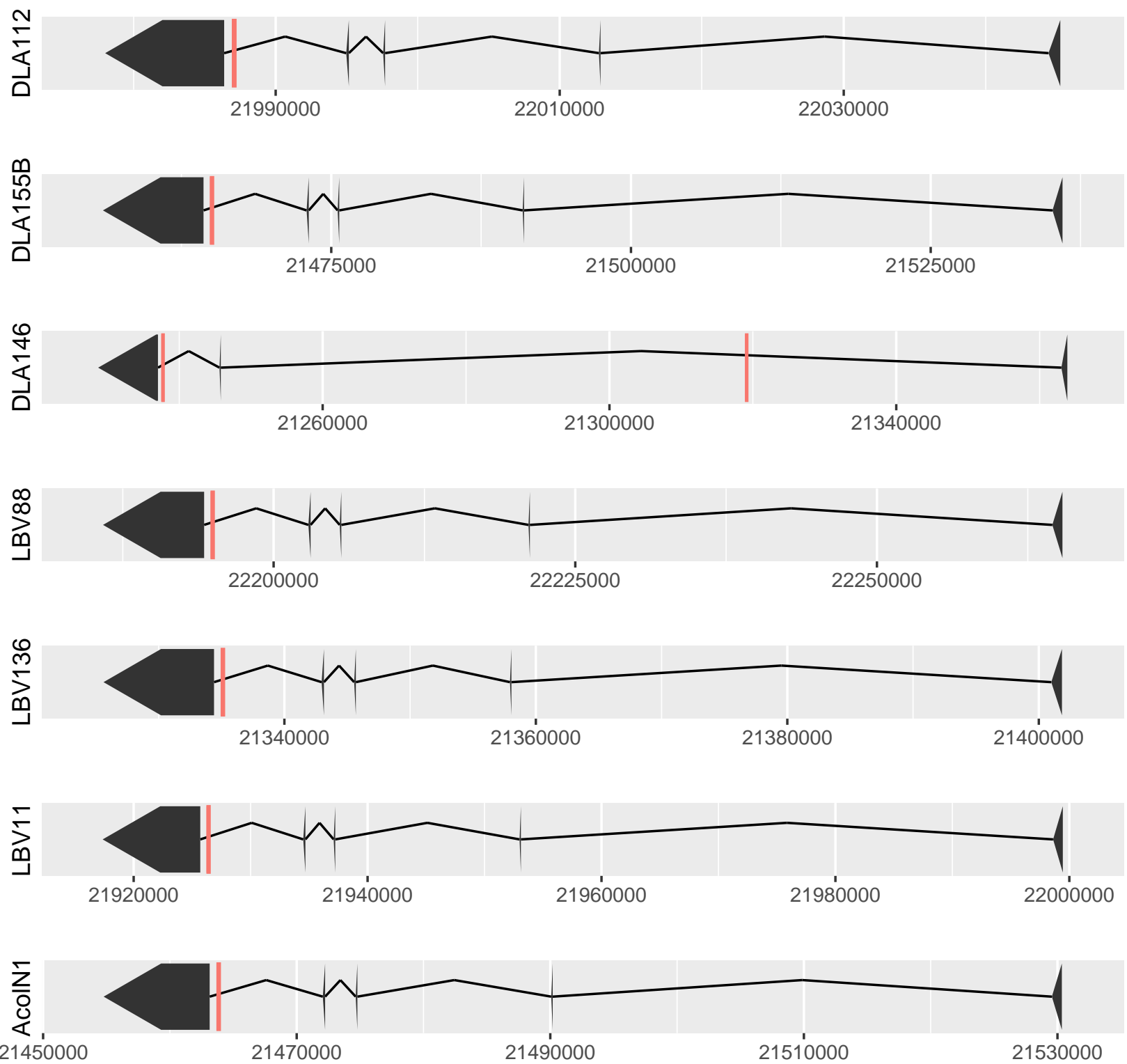
AGAP008190



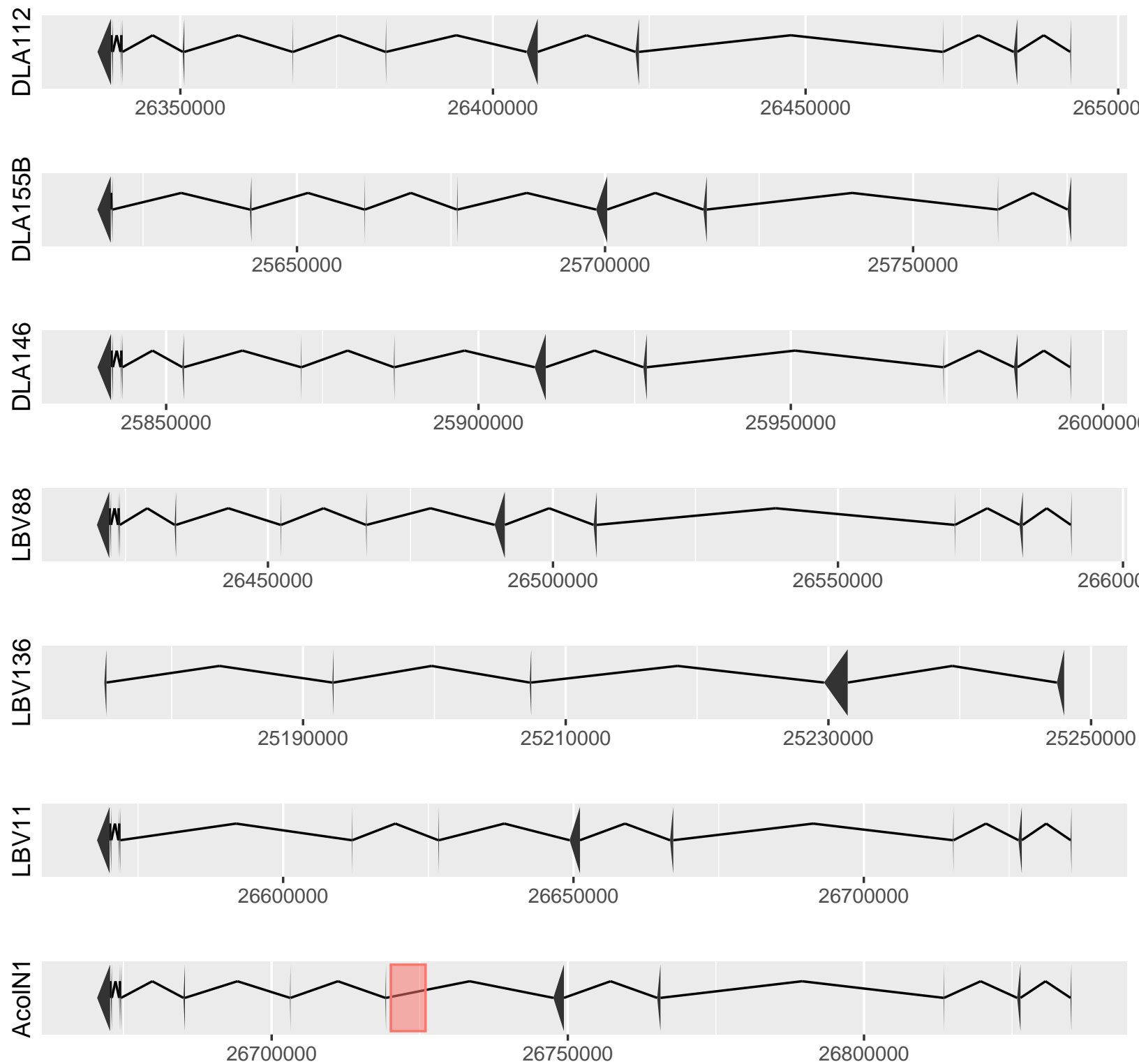
AGAP008535



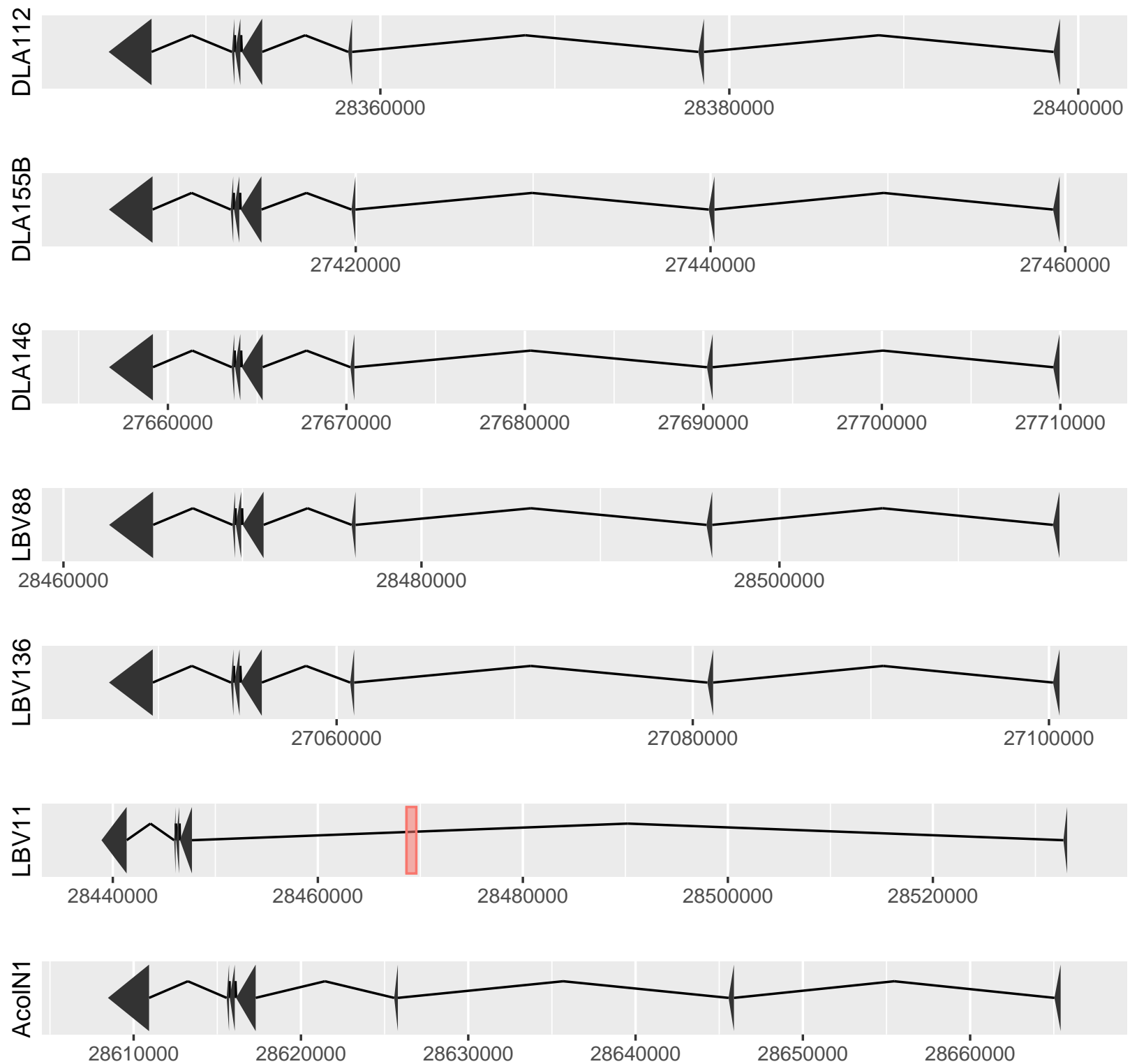
AGAP008944



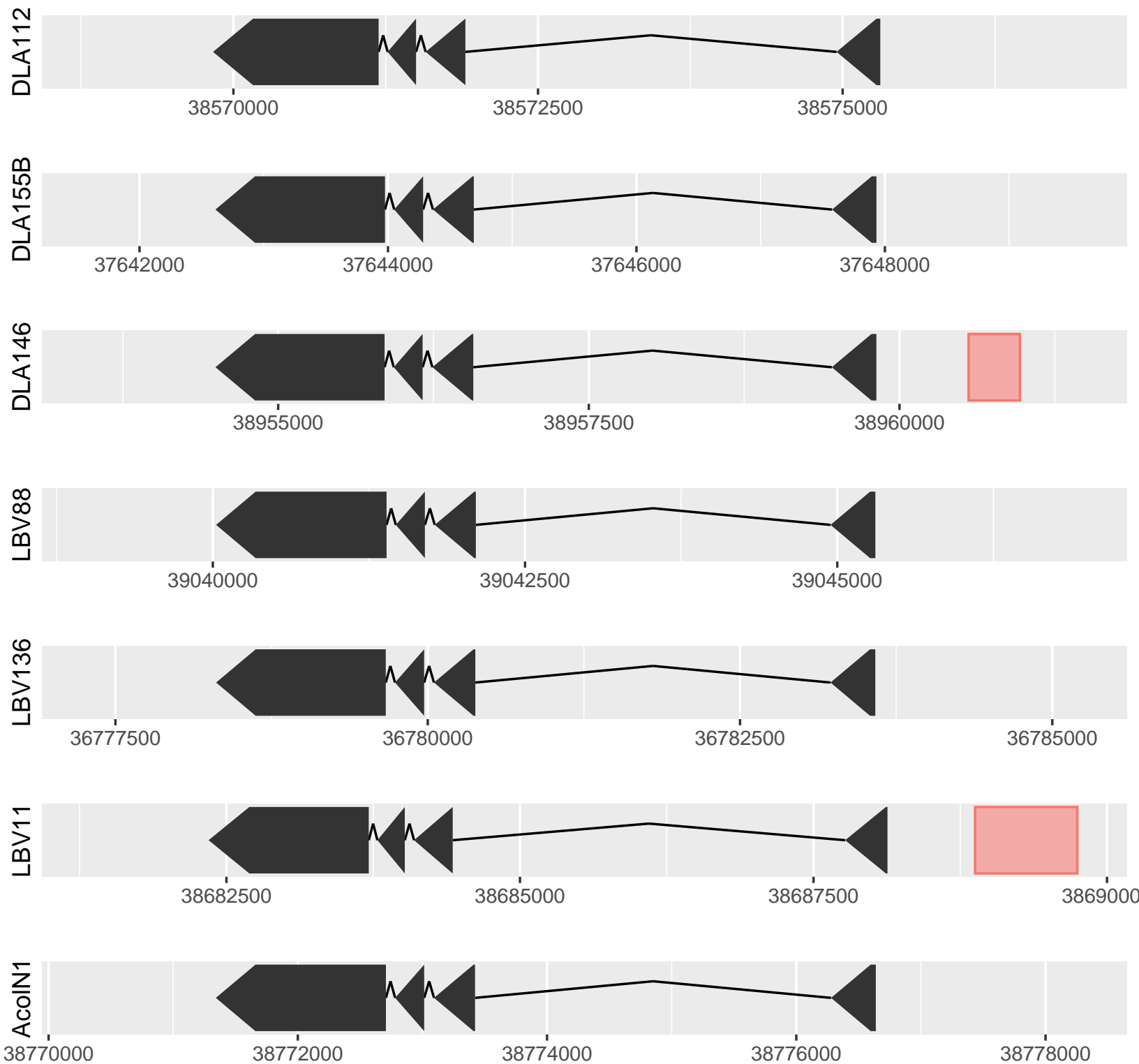
AGAP009112



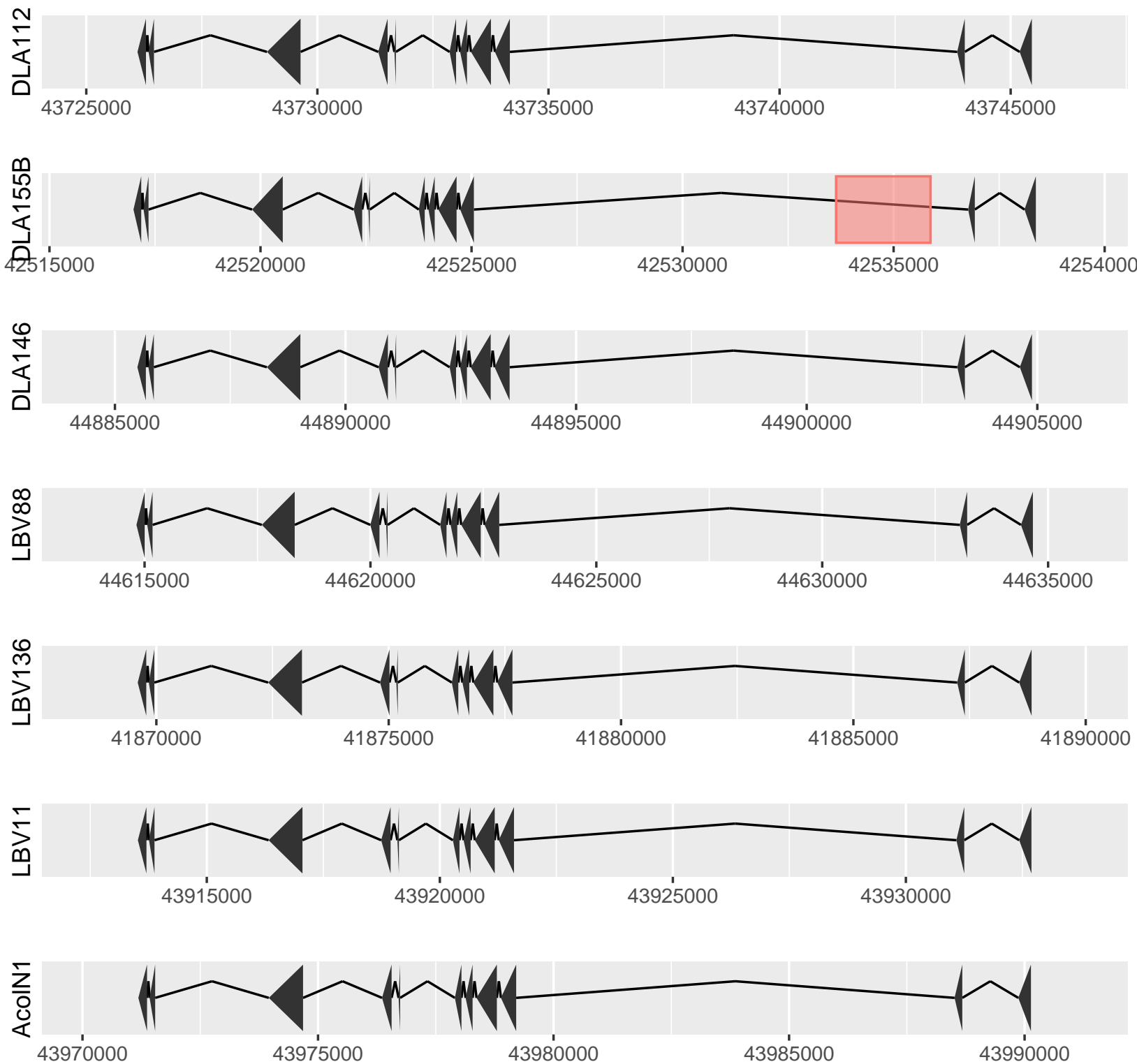
AGAP009158



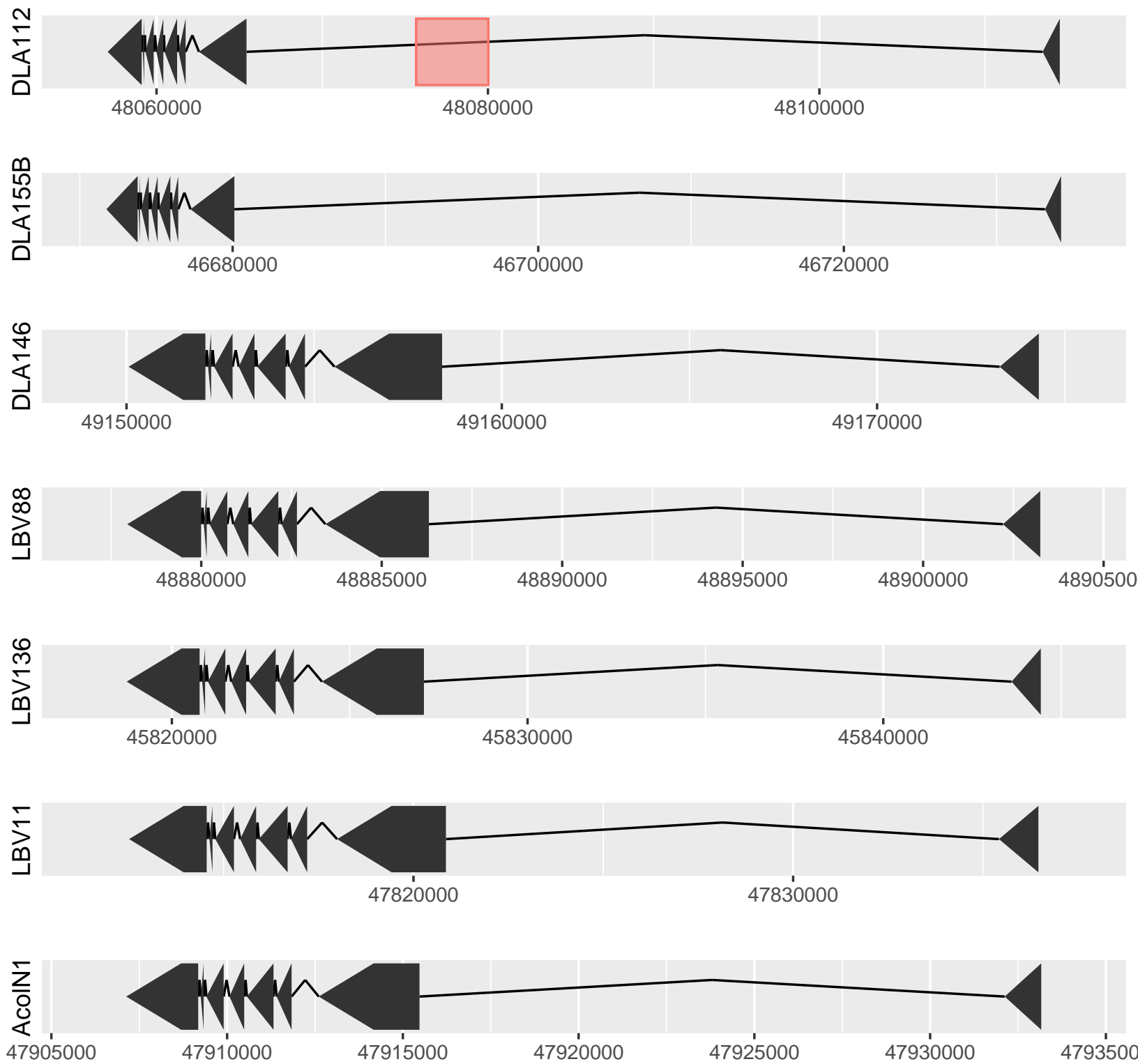
AGAP009668



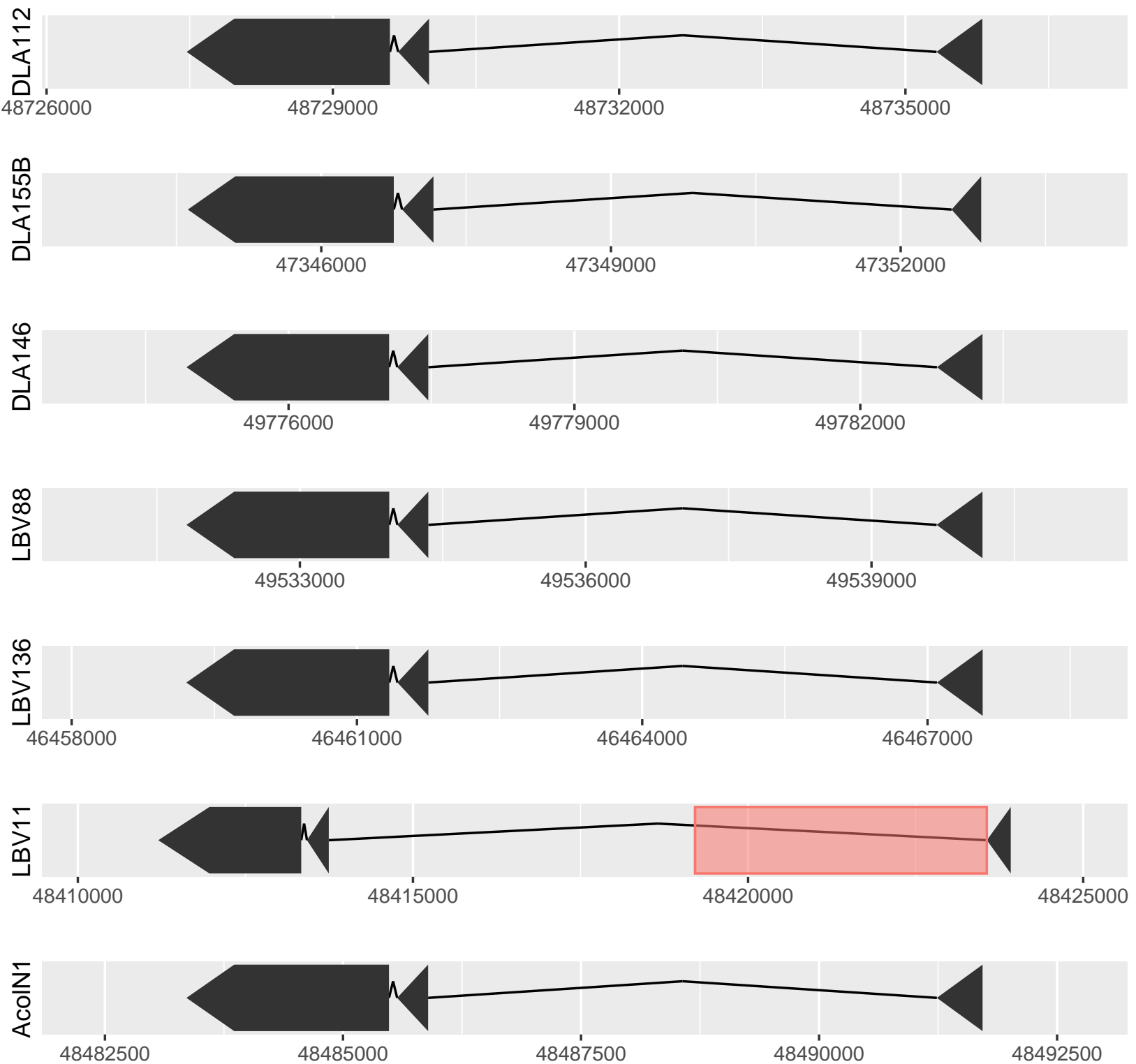
AGAP009774



AGAP009972

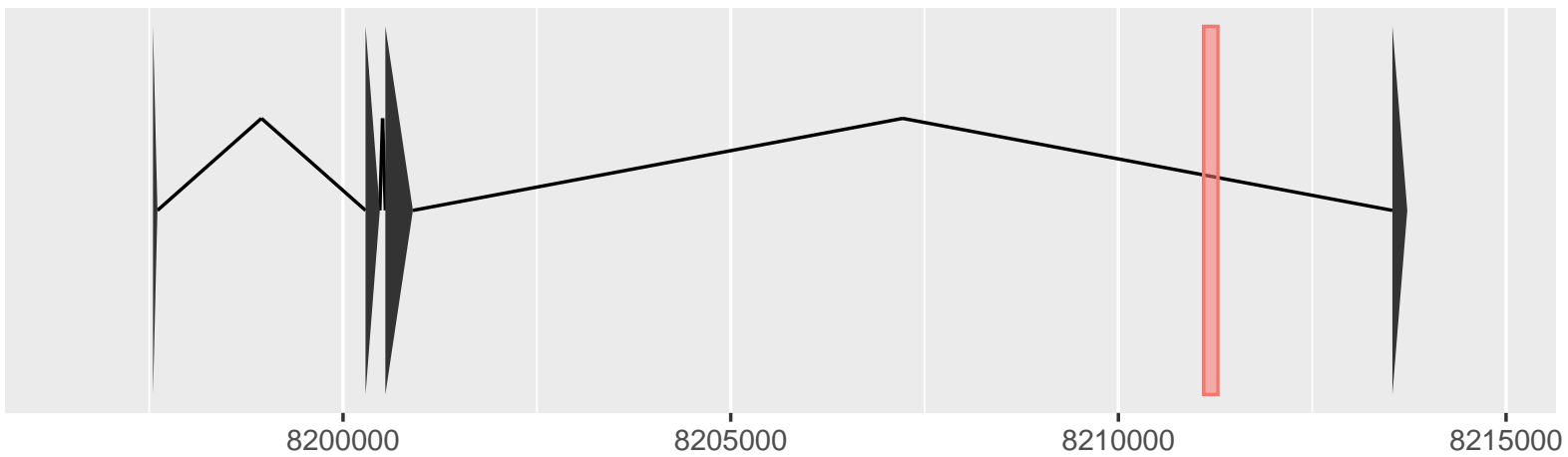


AGAP010008

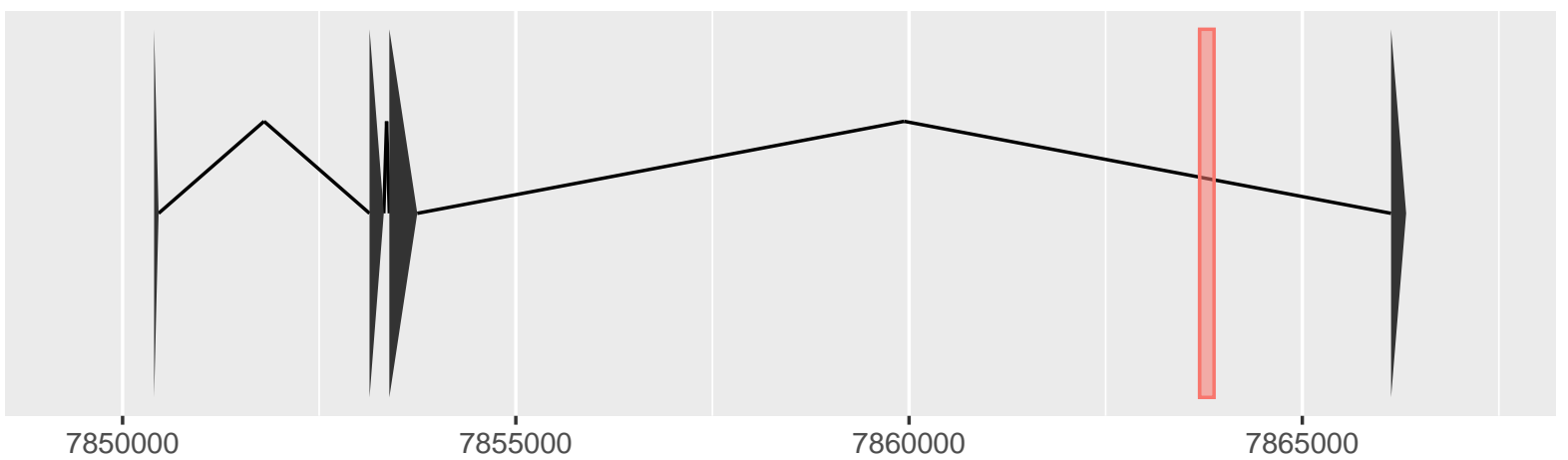


AGAP010620

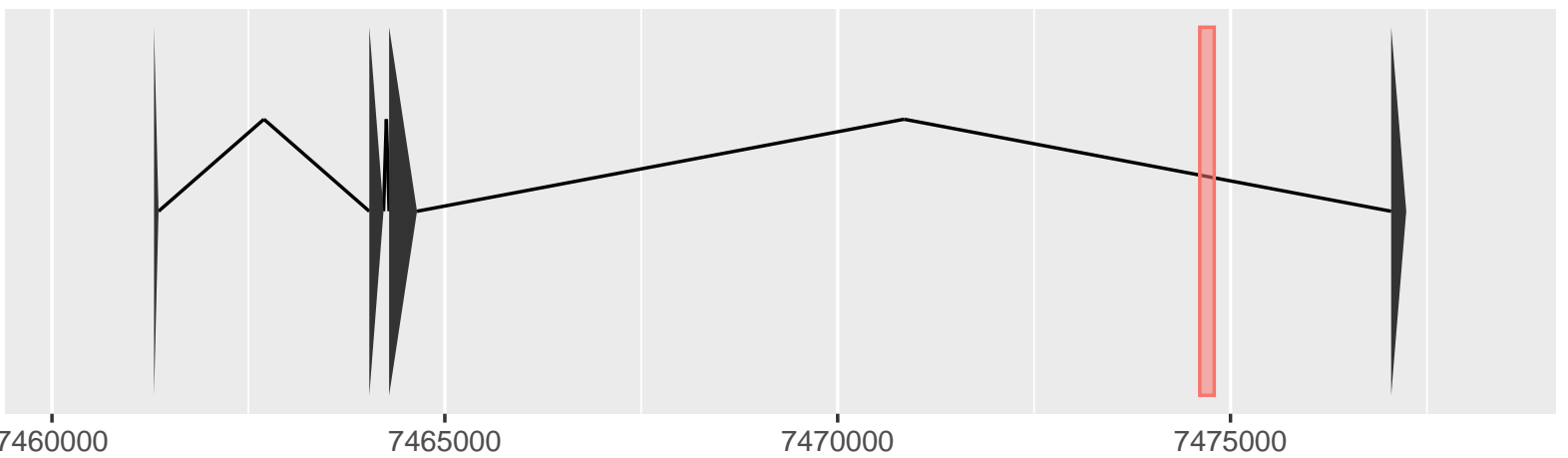
DLA146



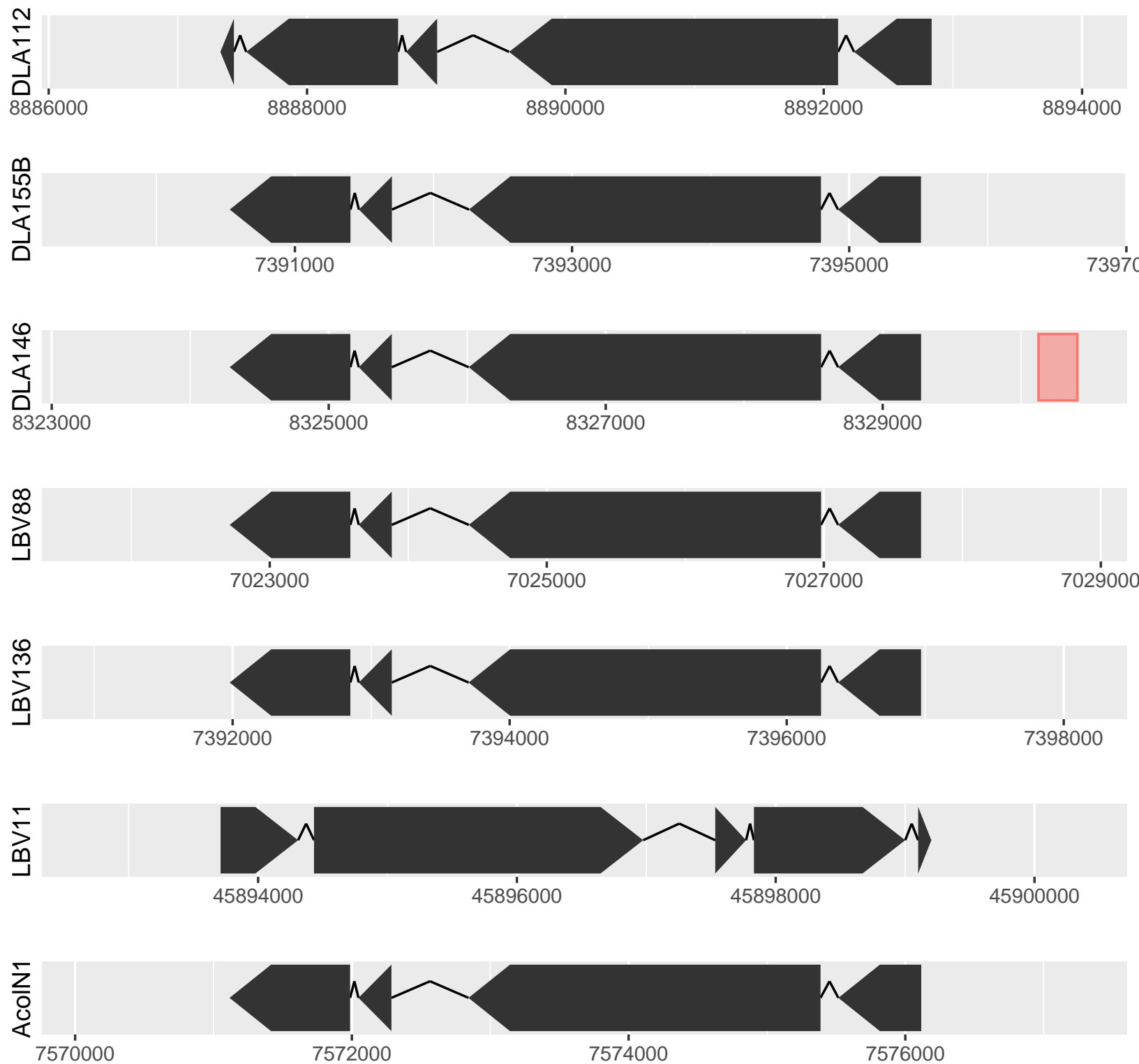
LBV11



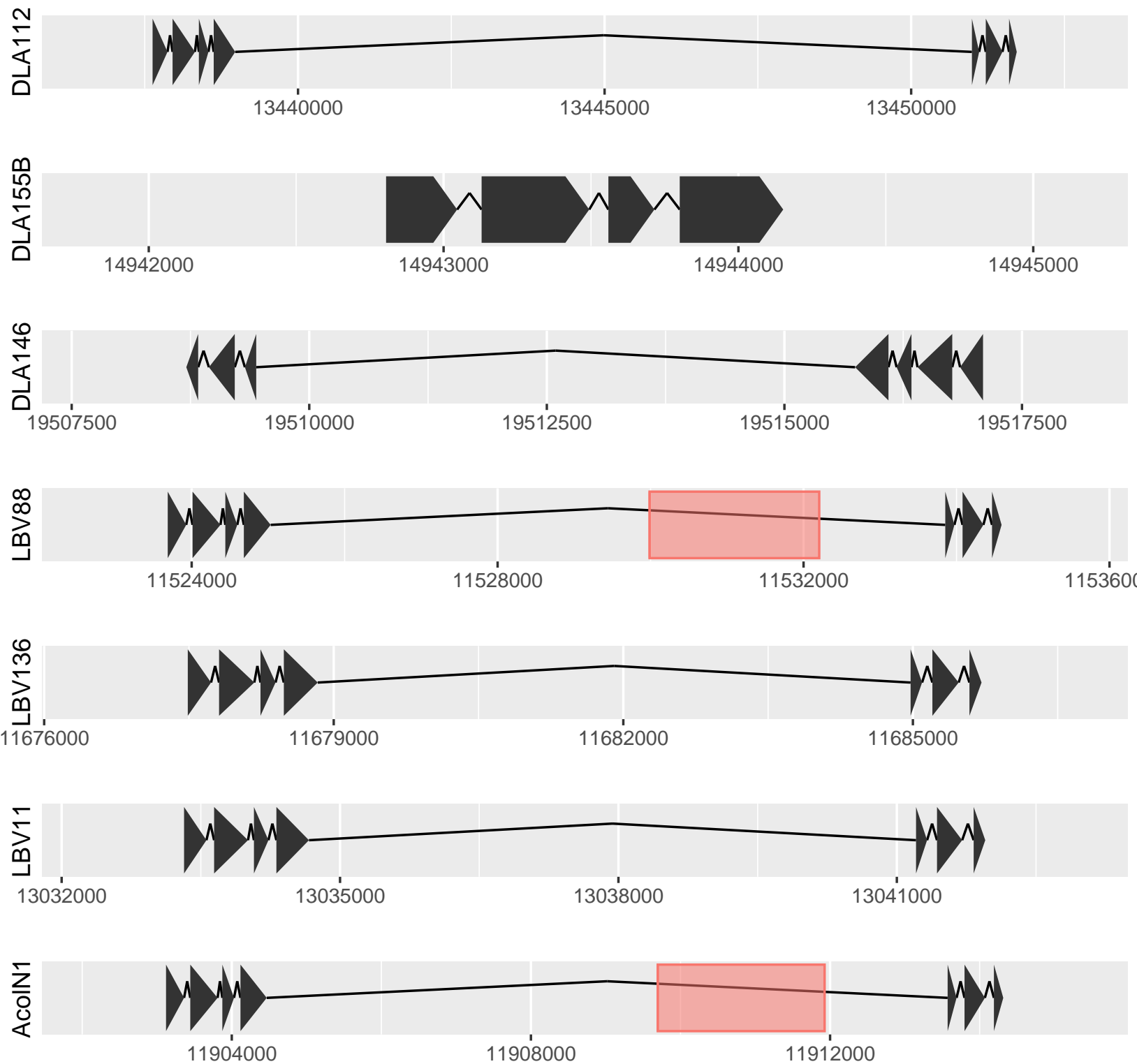
AcoIN1



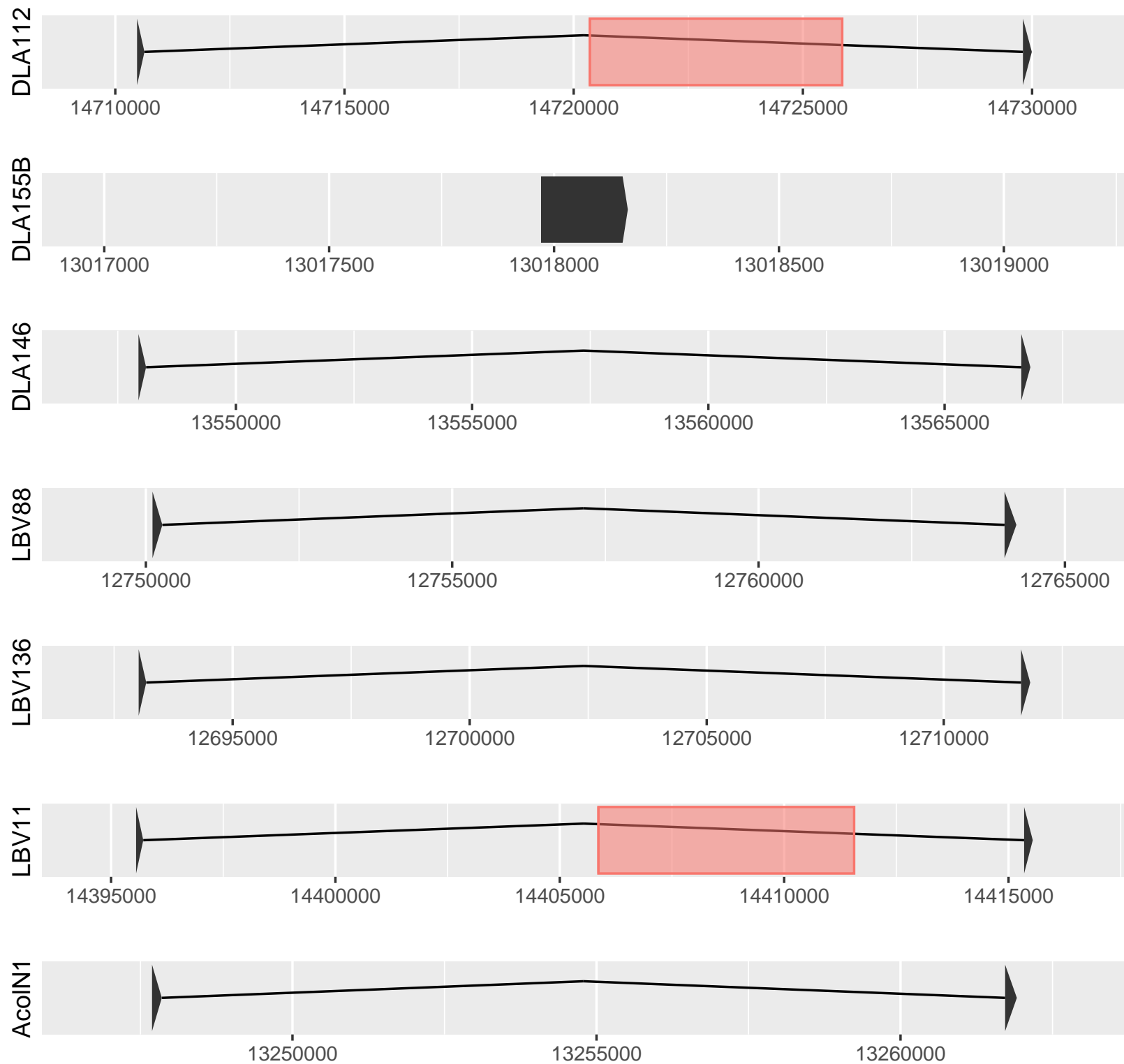
AGAP010626



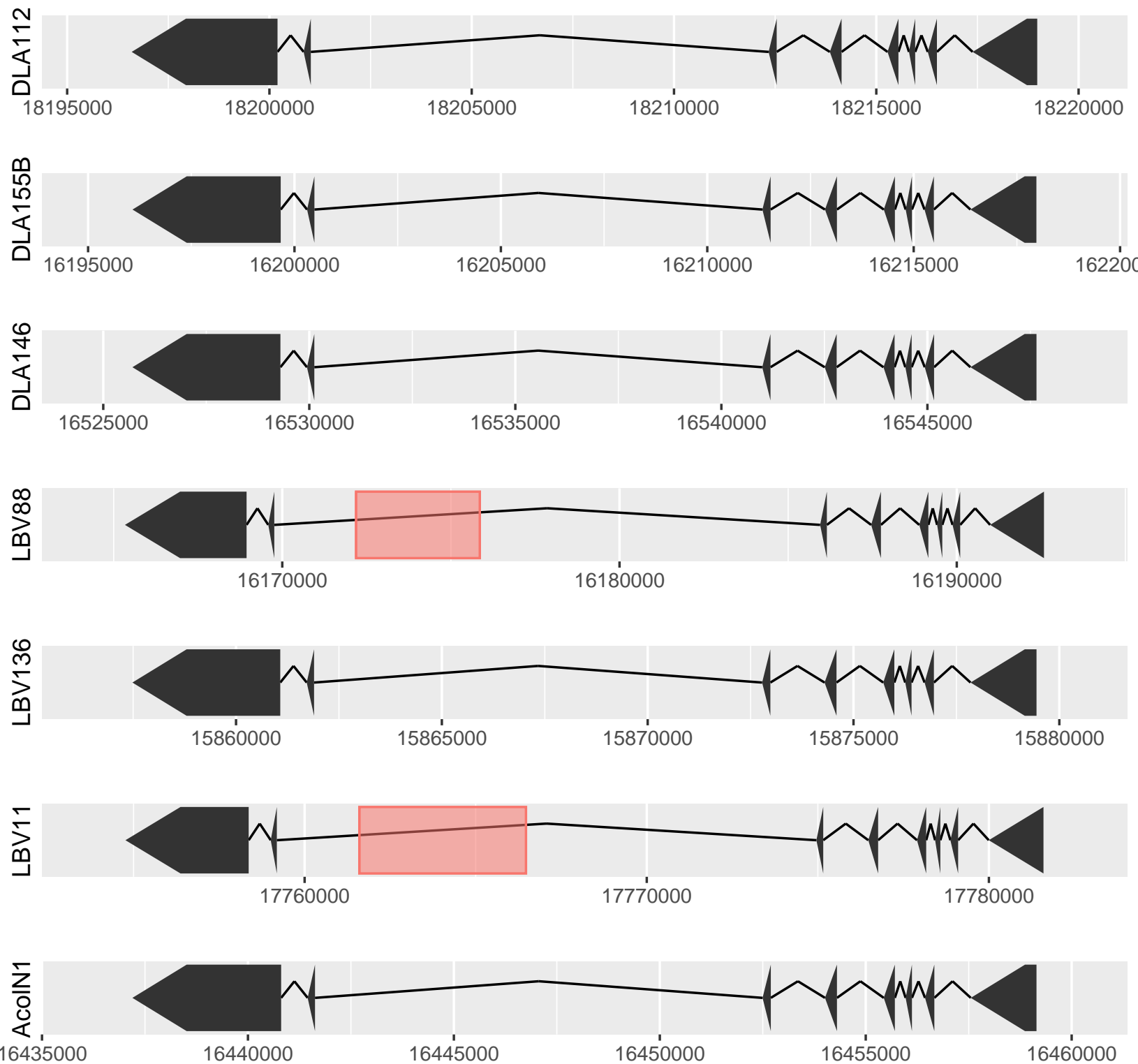
AGAP010850



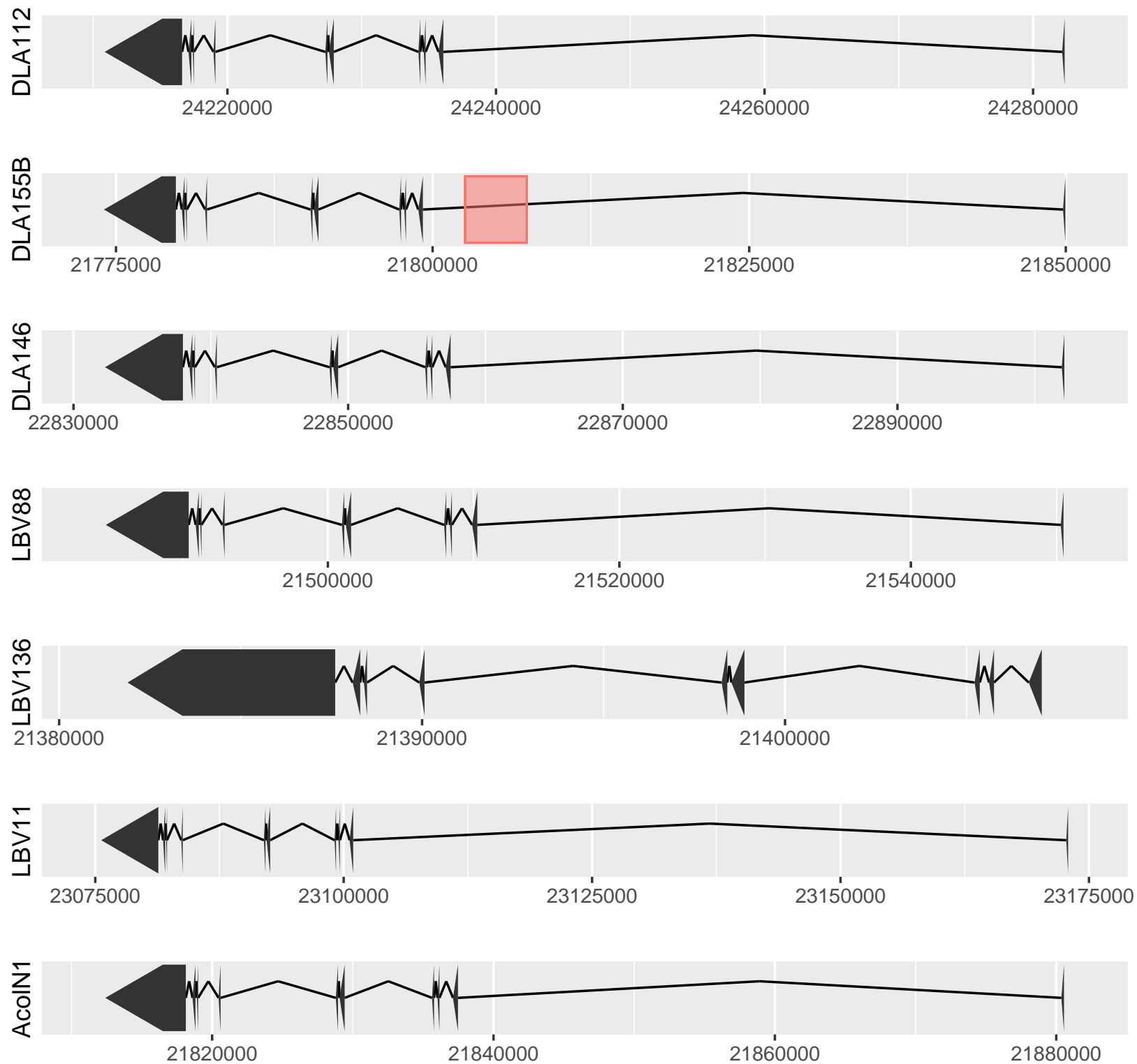
AGAP010922



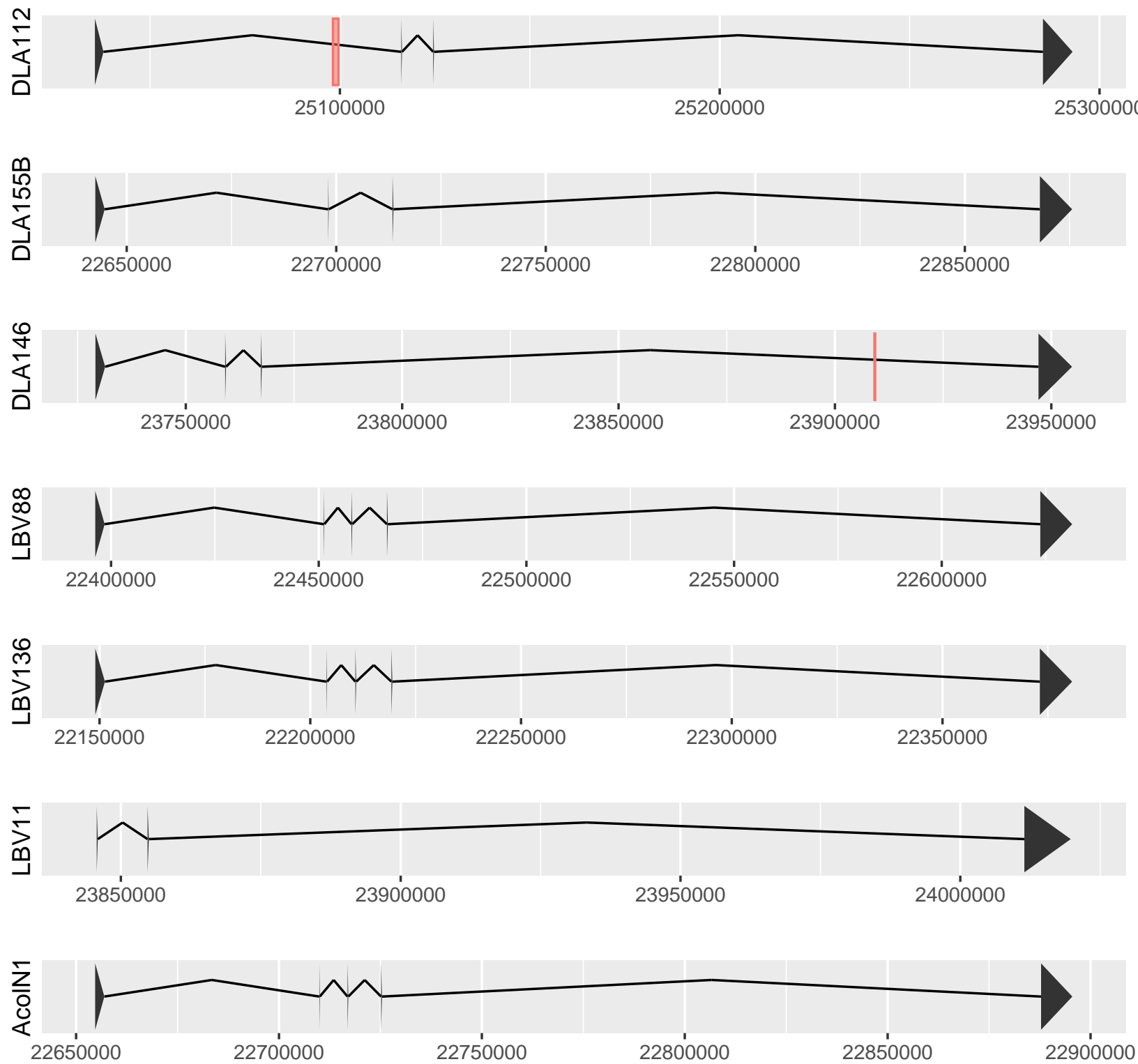
AGAP011094



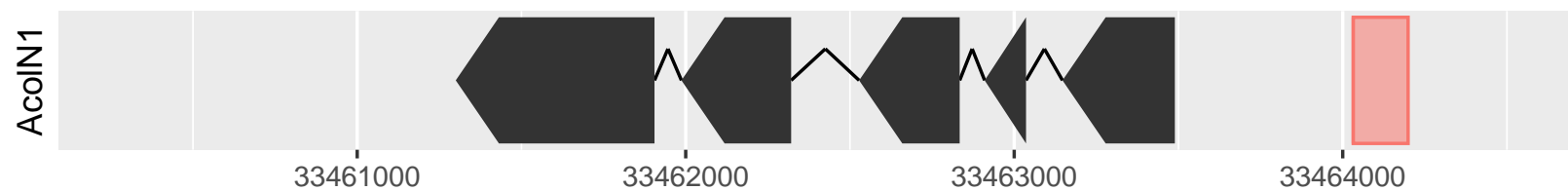
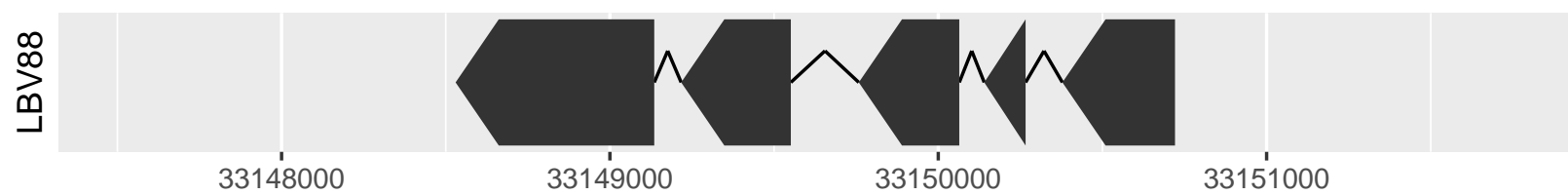
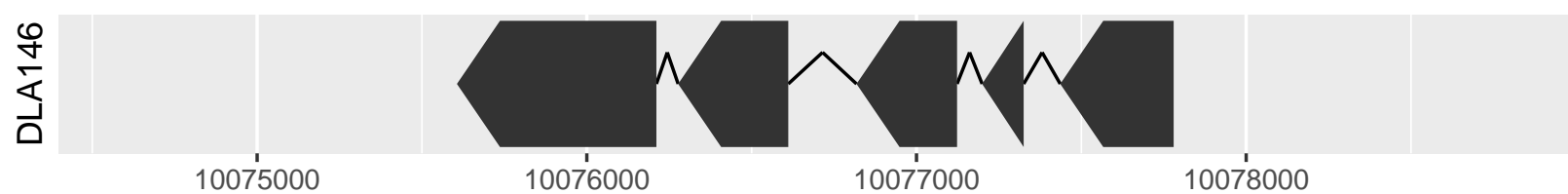
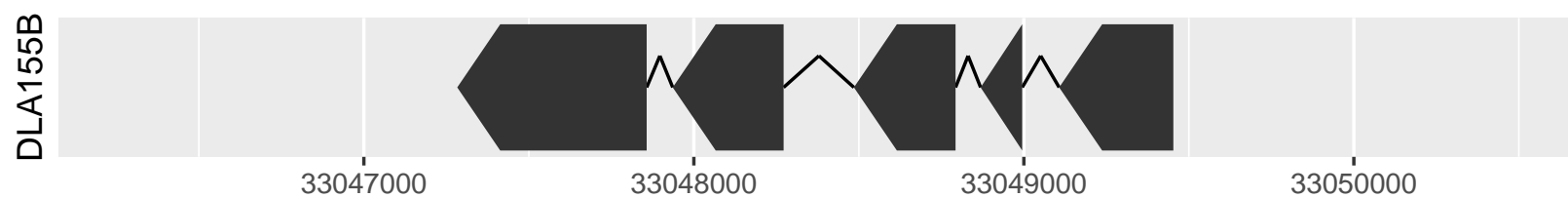
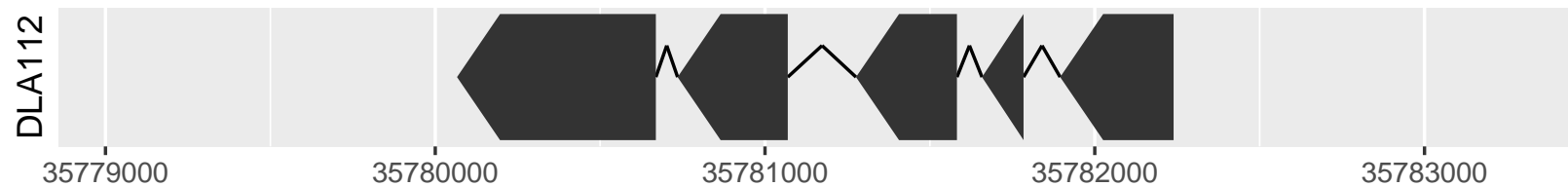
AGAP011360



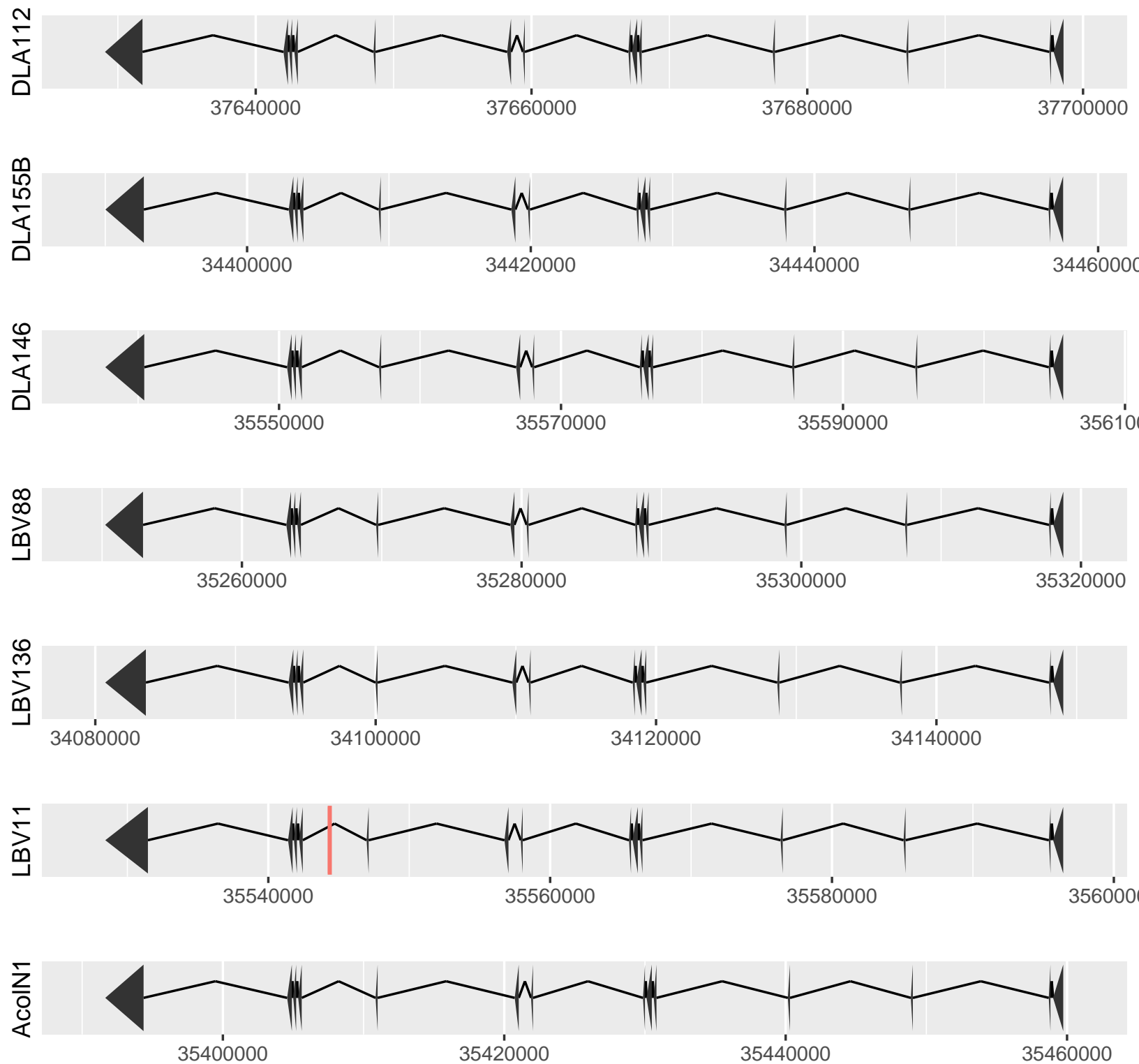
AGAP011379



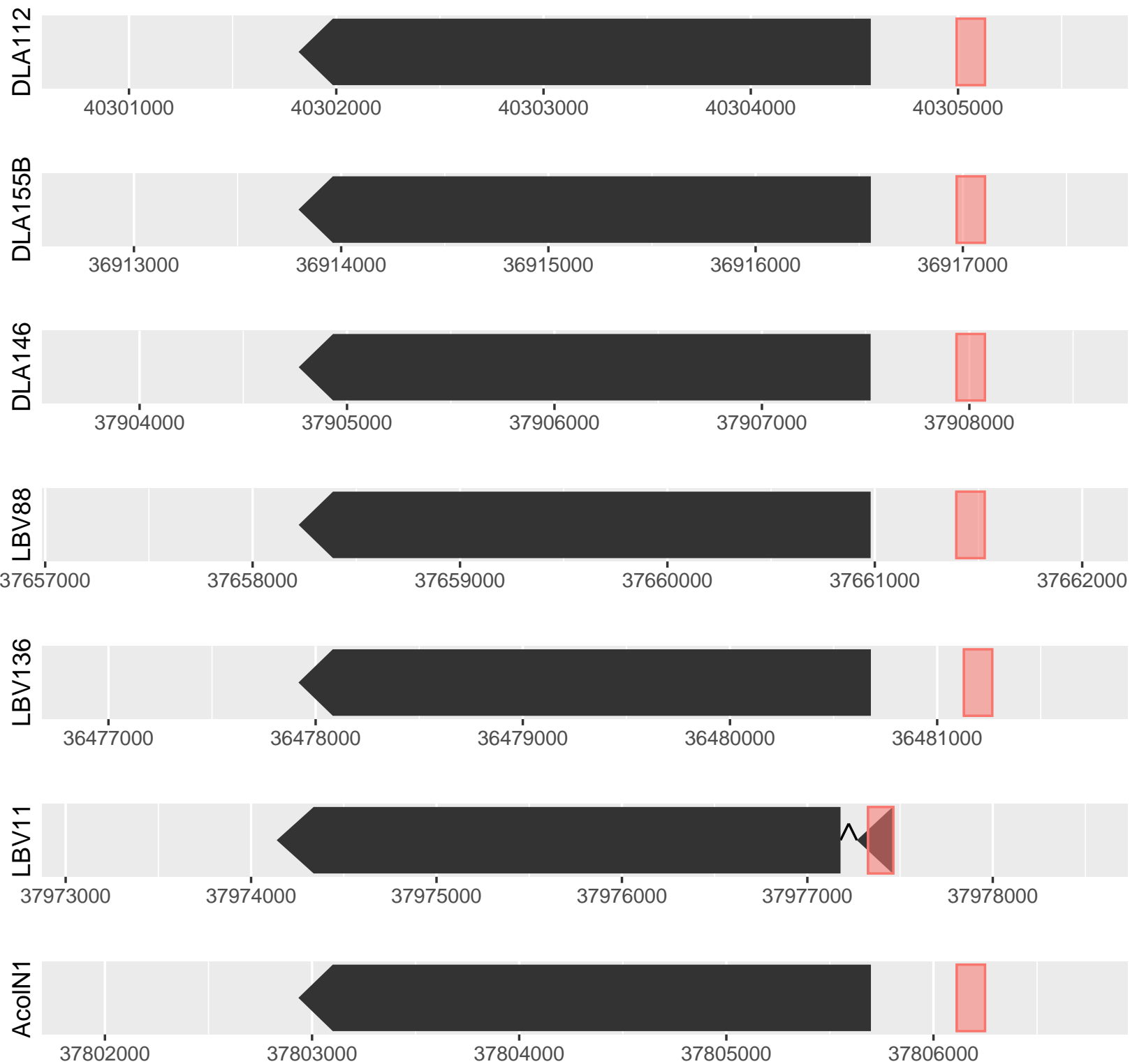
AGAP011794



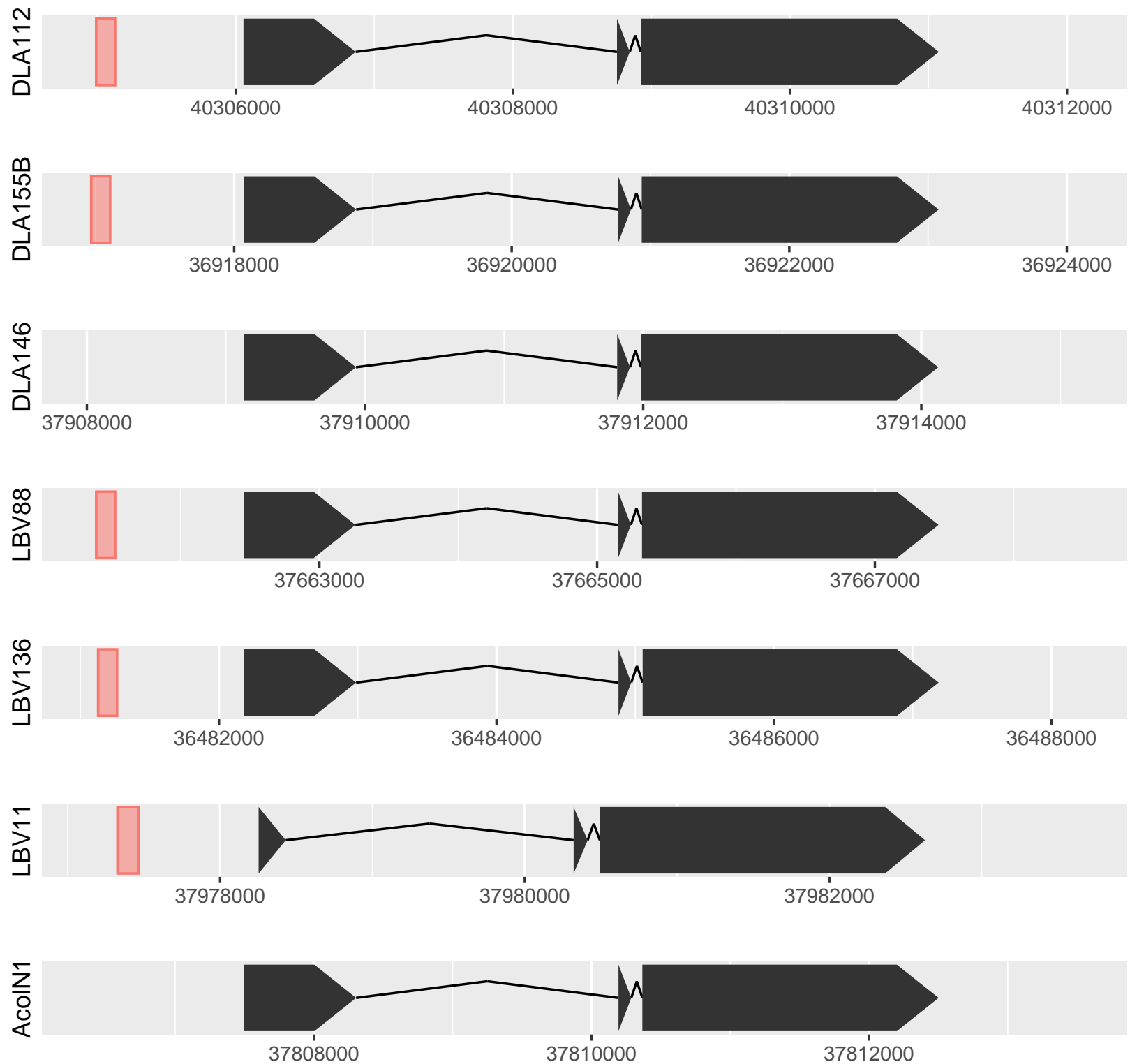
AGAP011916



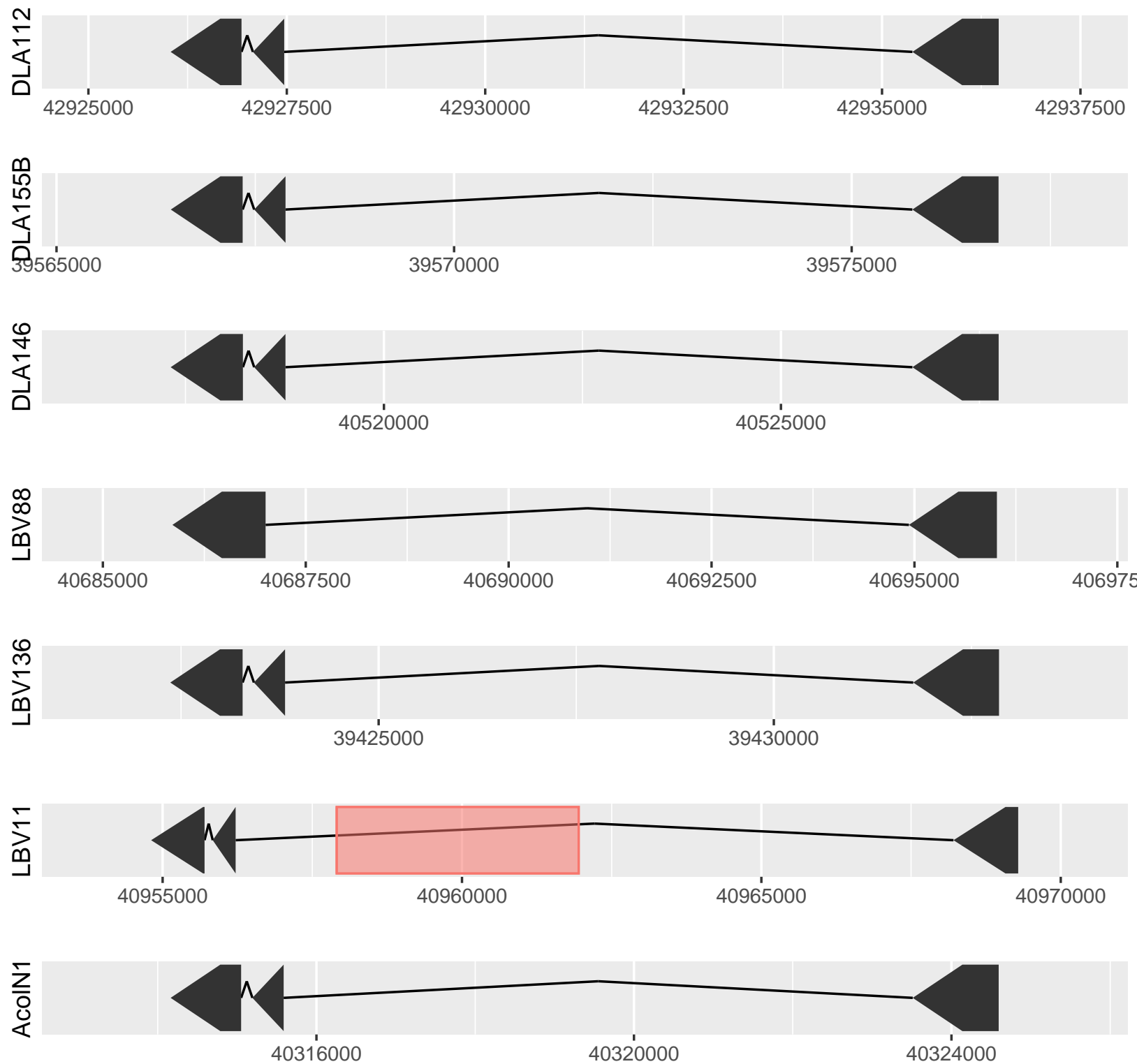
AGAP012092



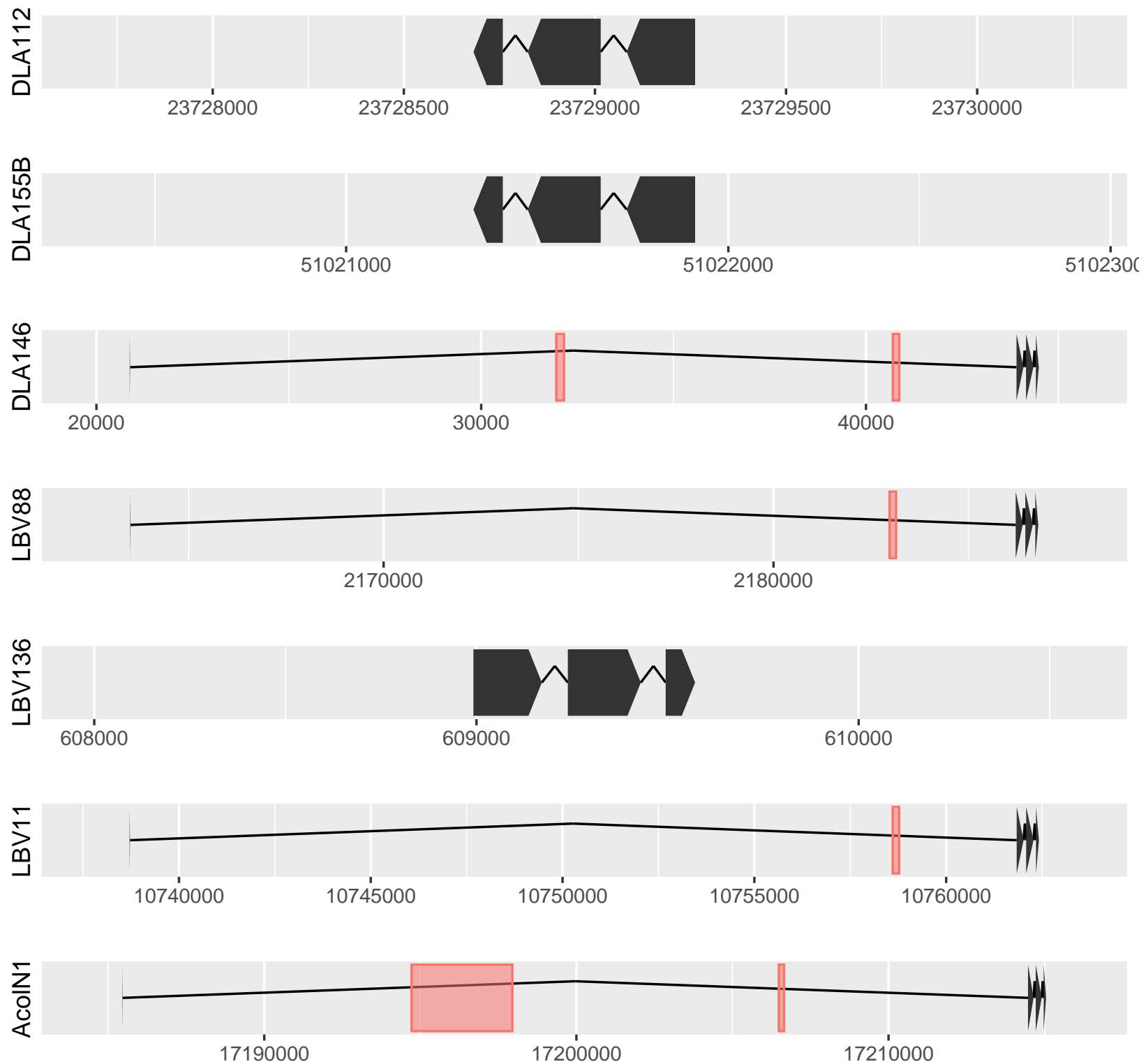
AGAP012093



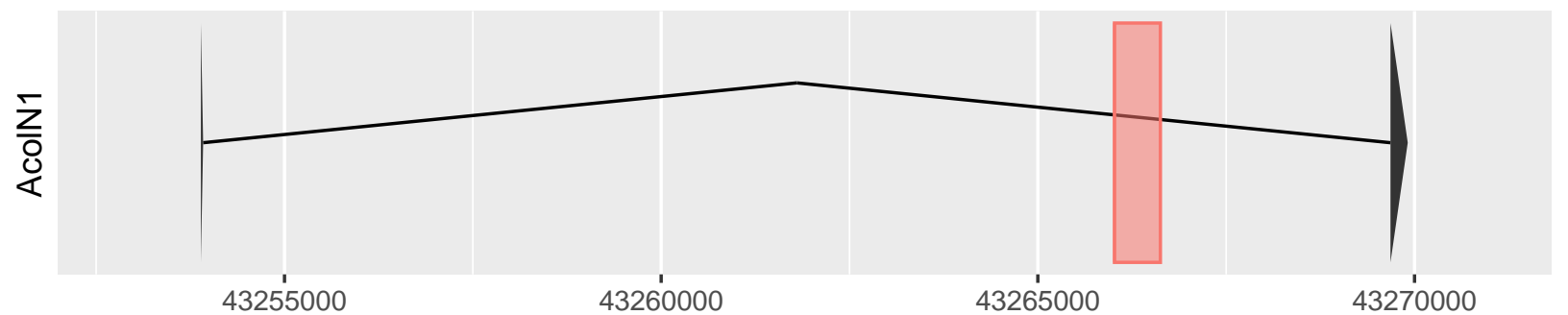
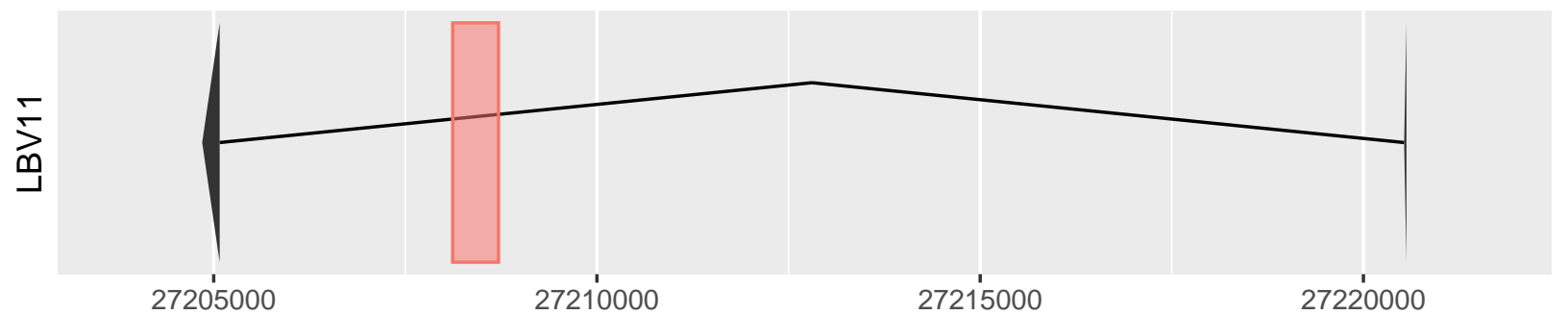
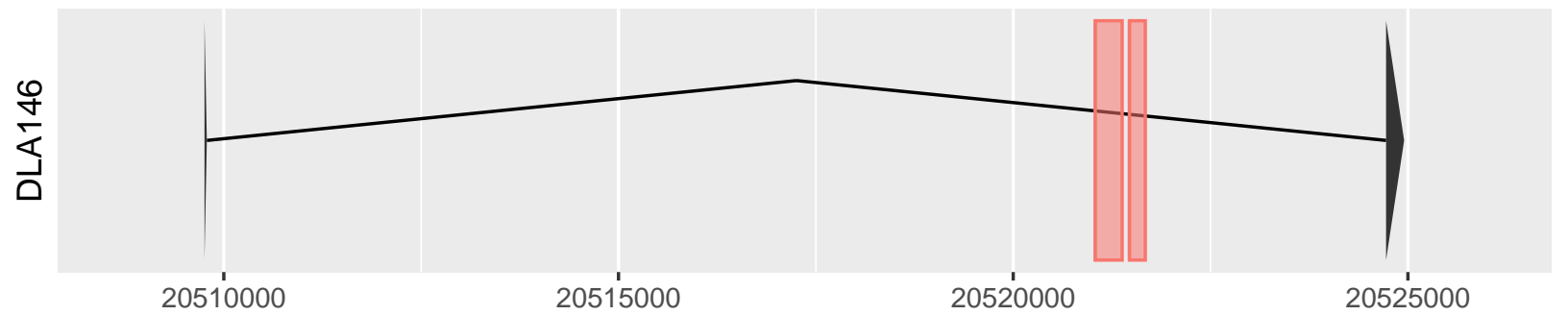
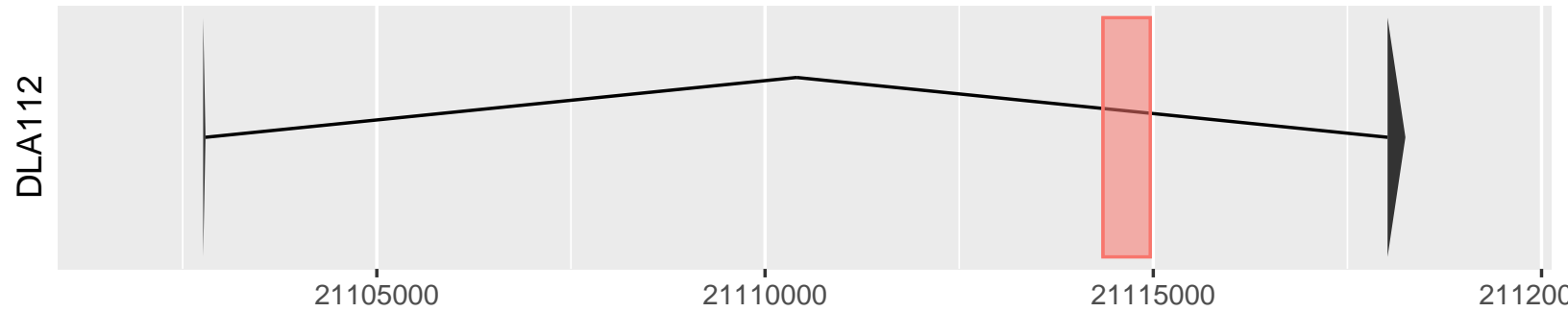
AGAP012316



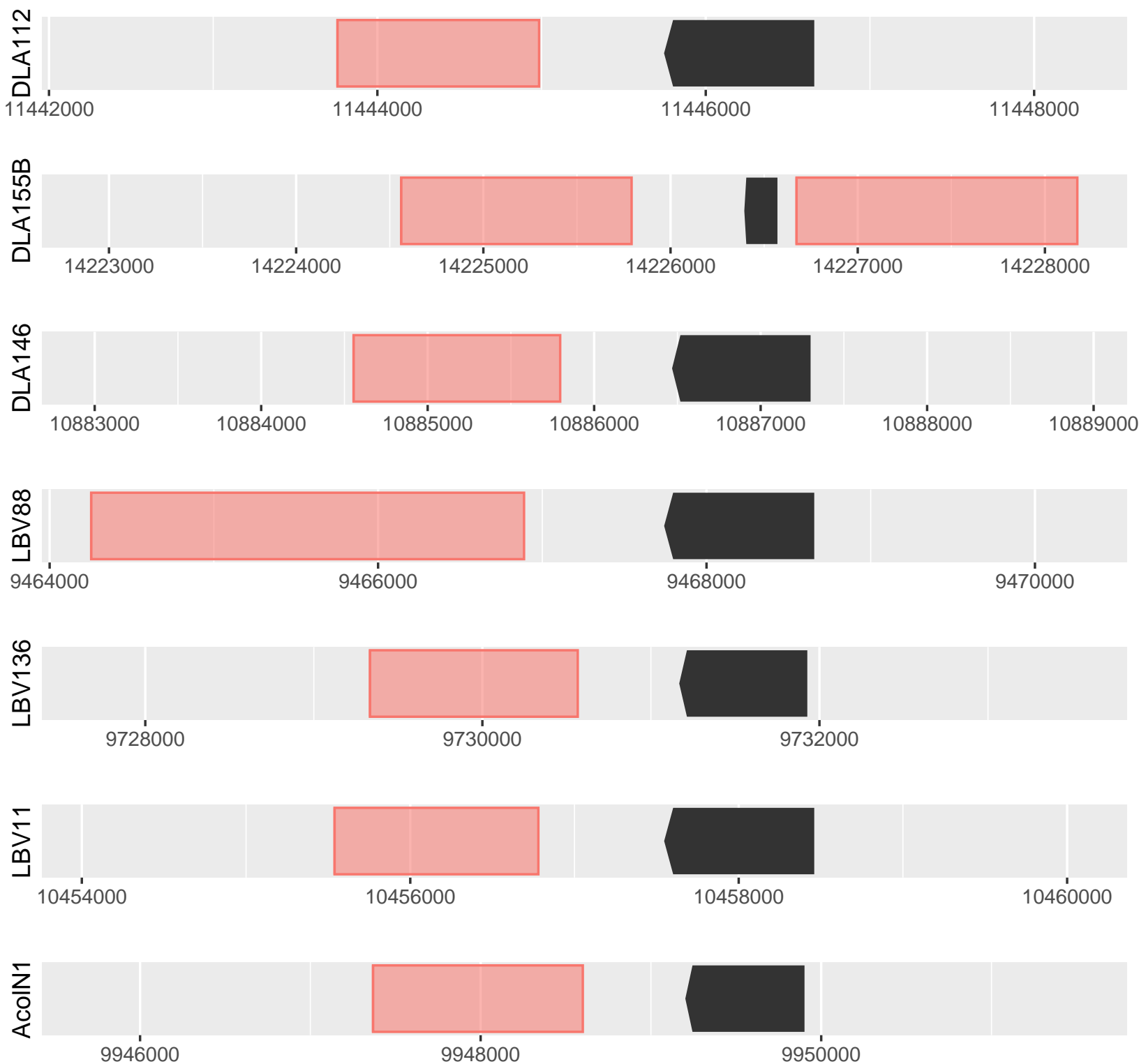
AGAP012466



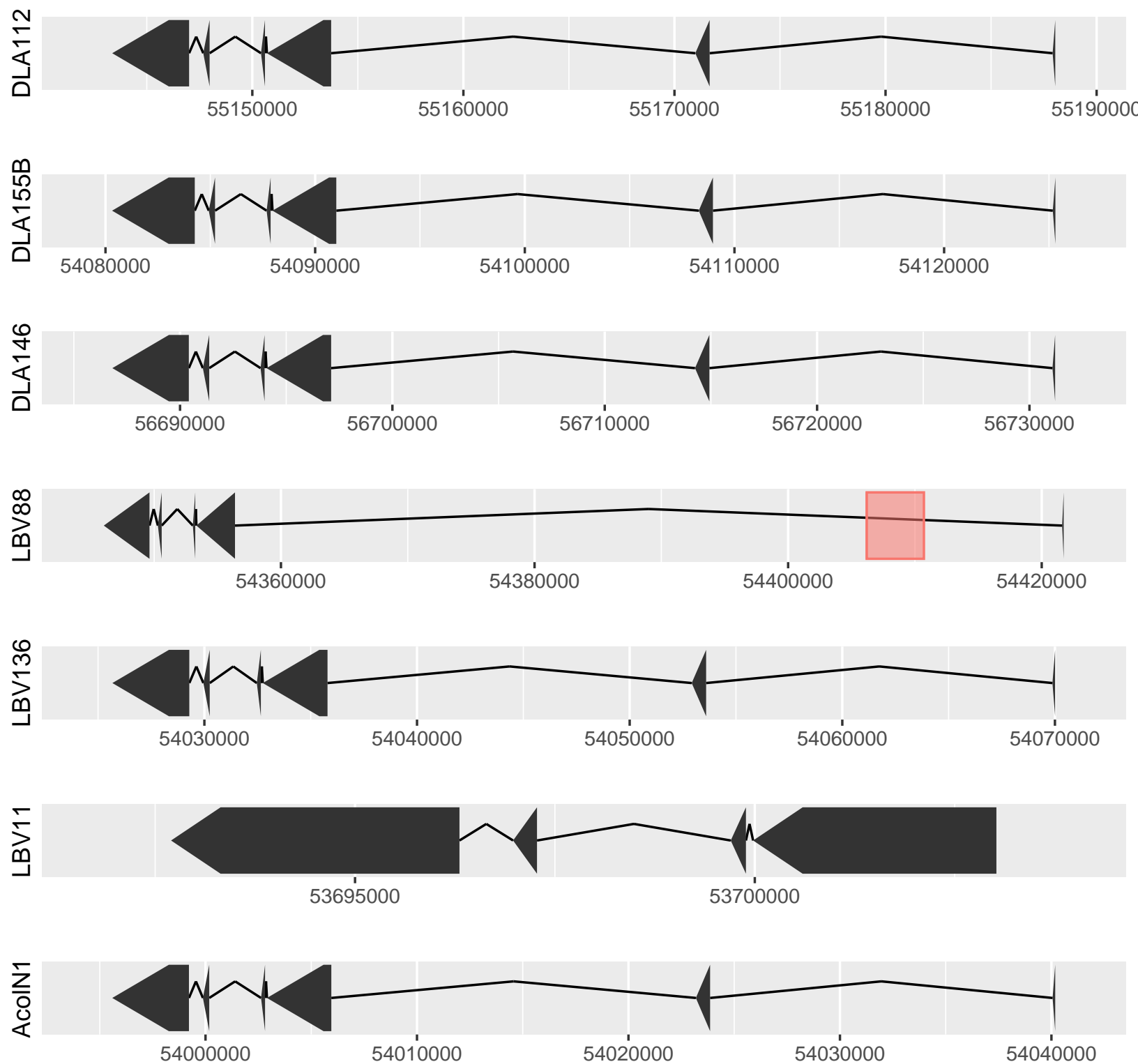
AGAP012484



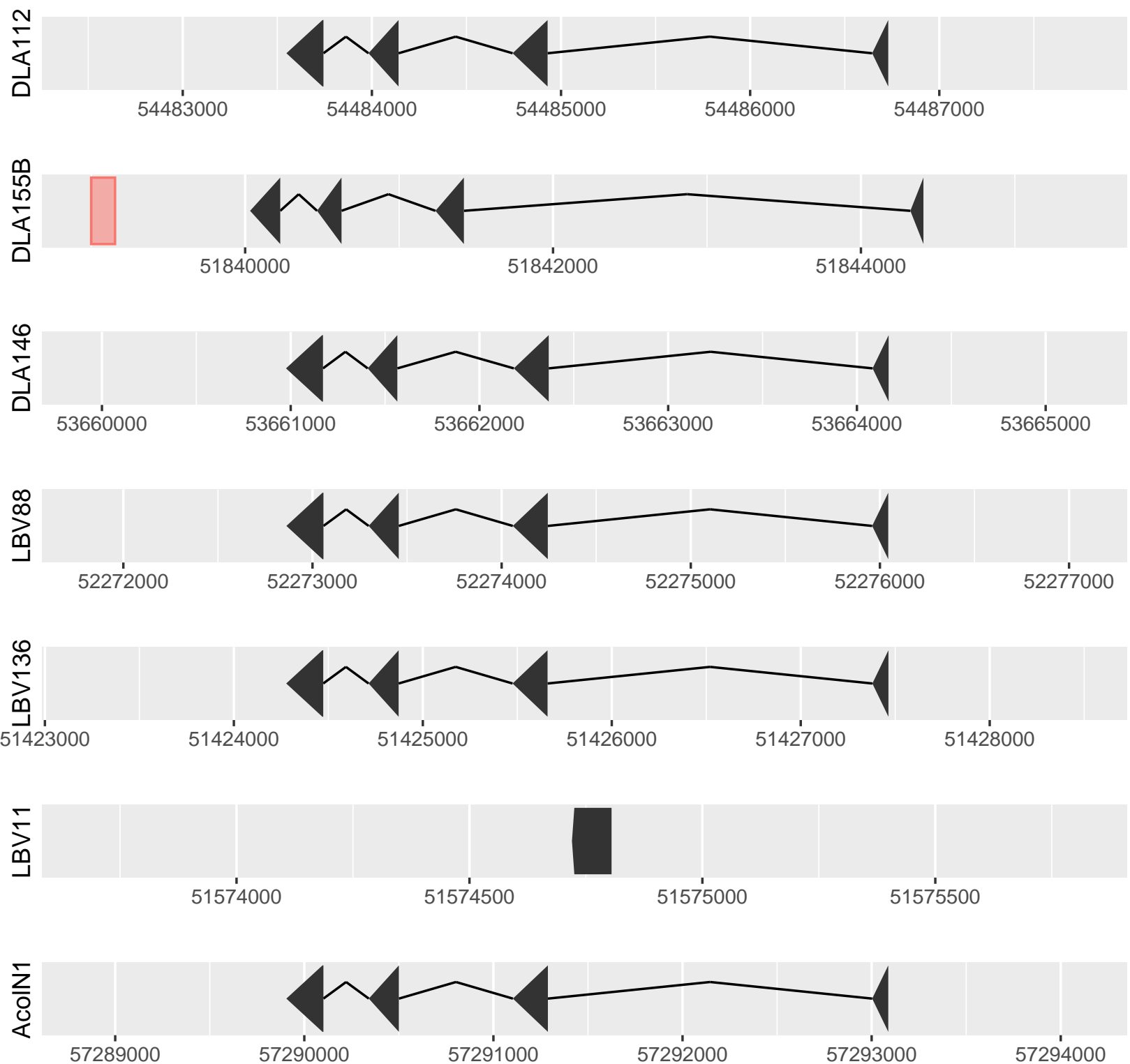
AGAP012916



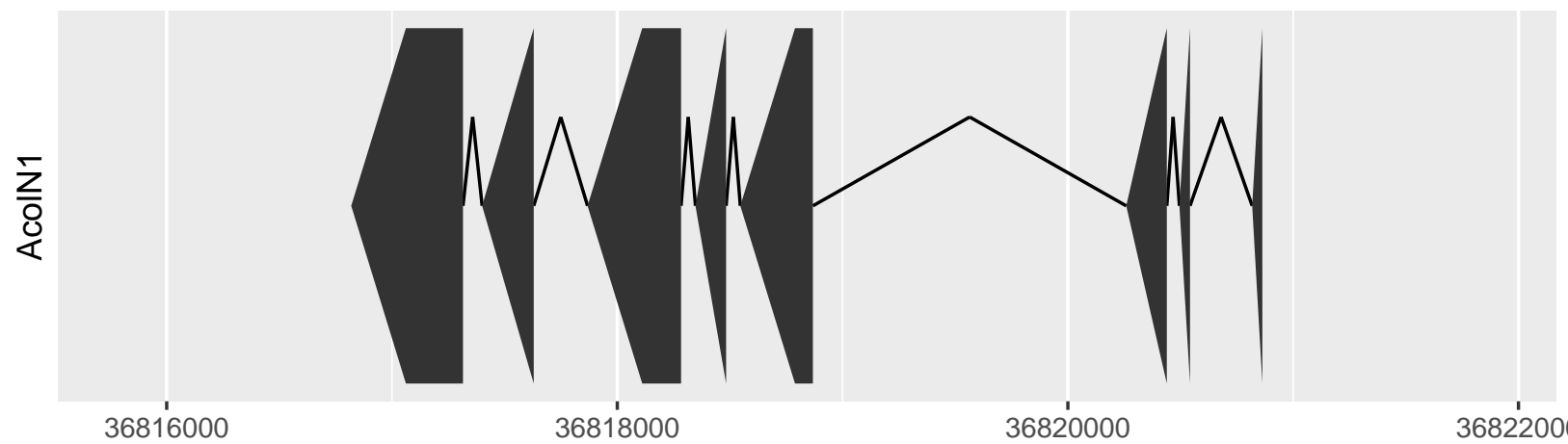
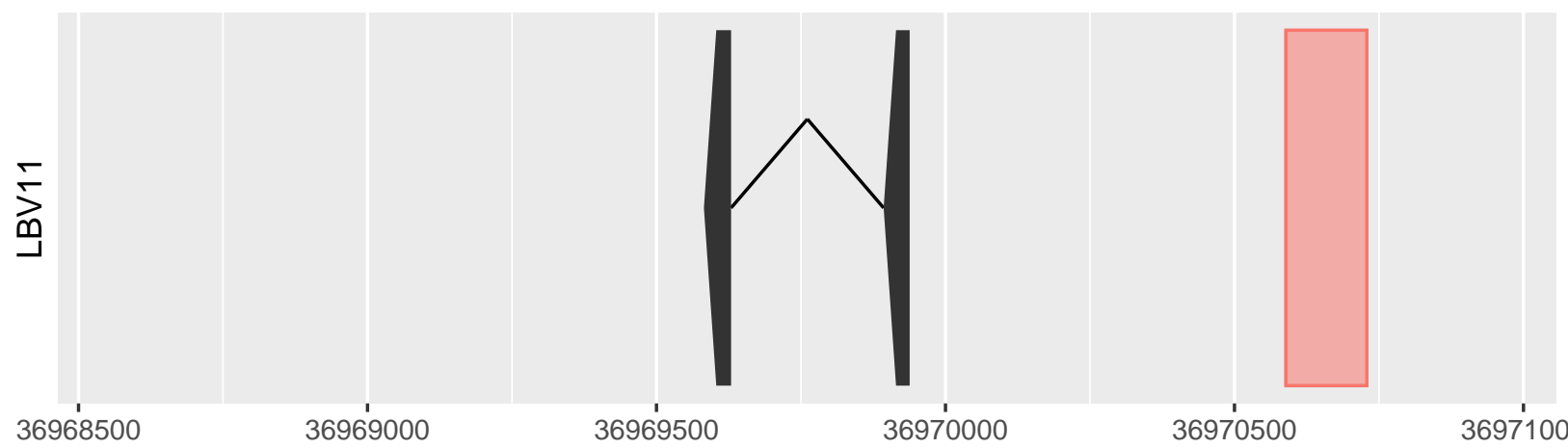
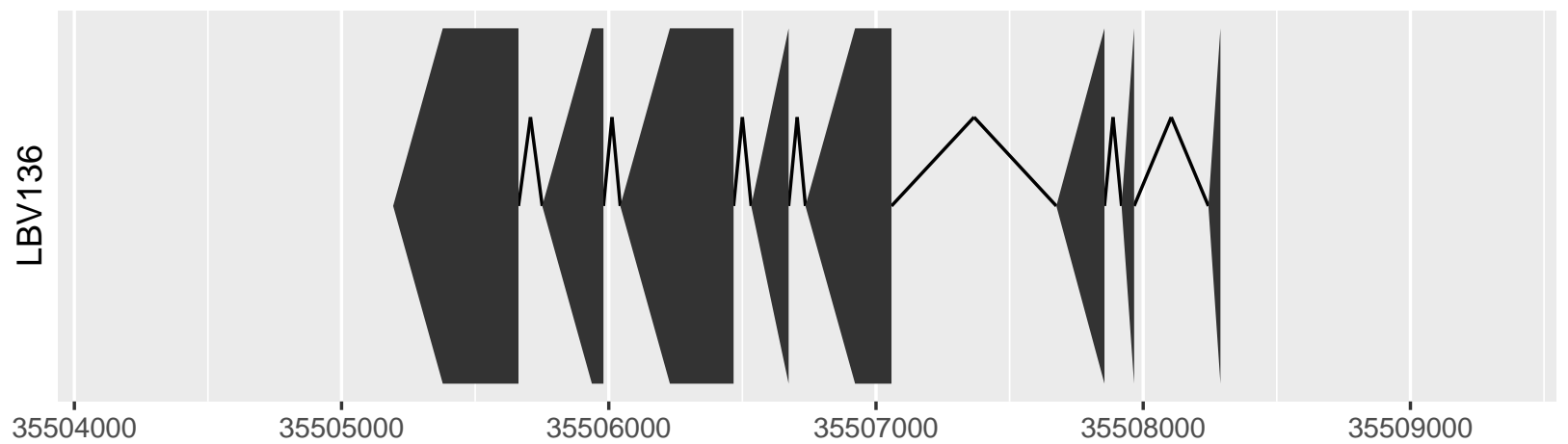
AGAP013057



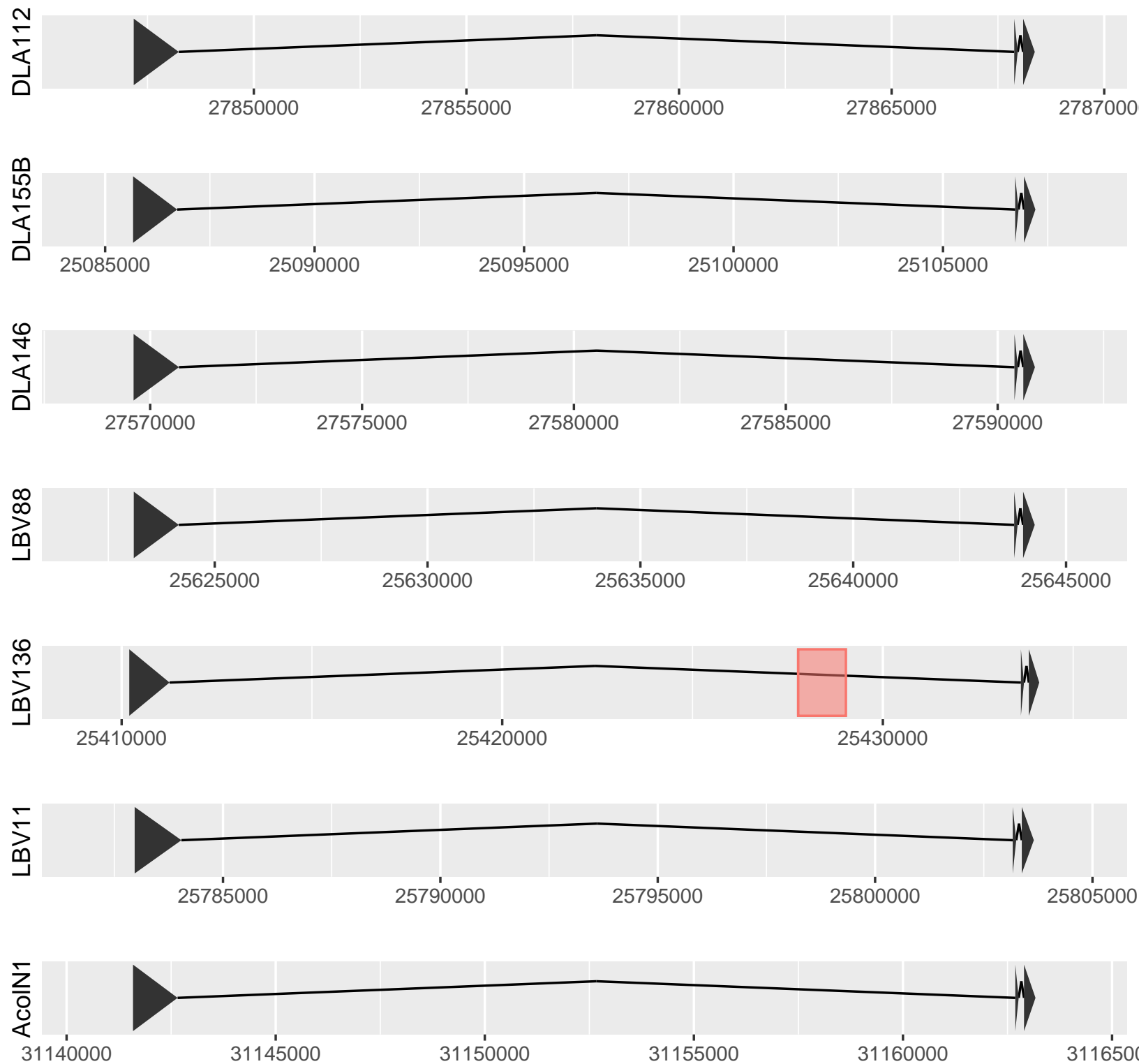
AGAP013228



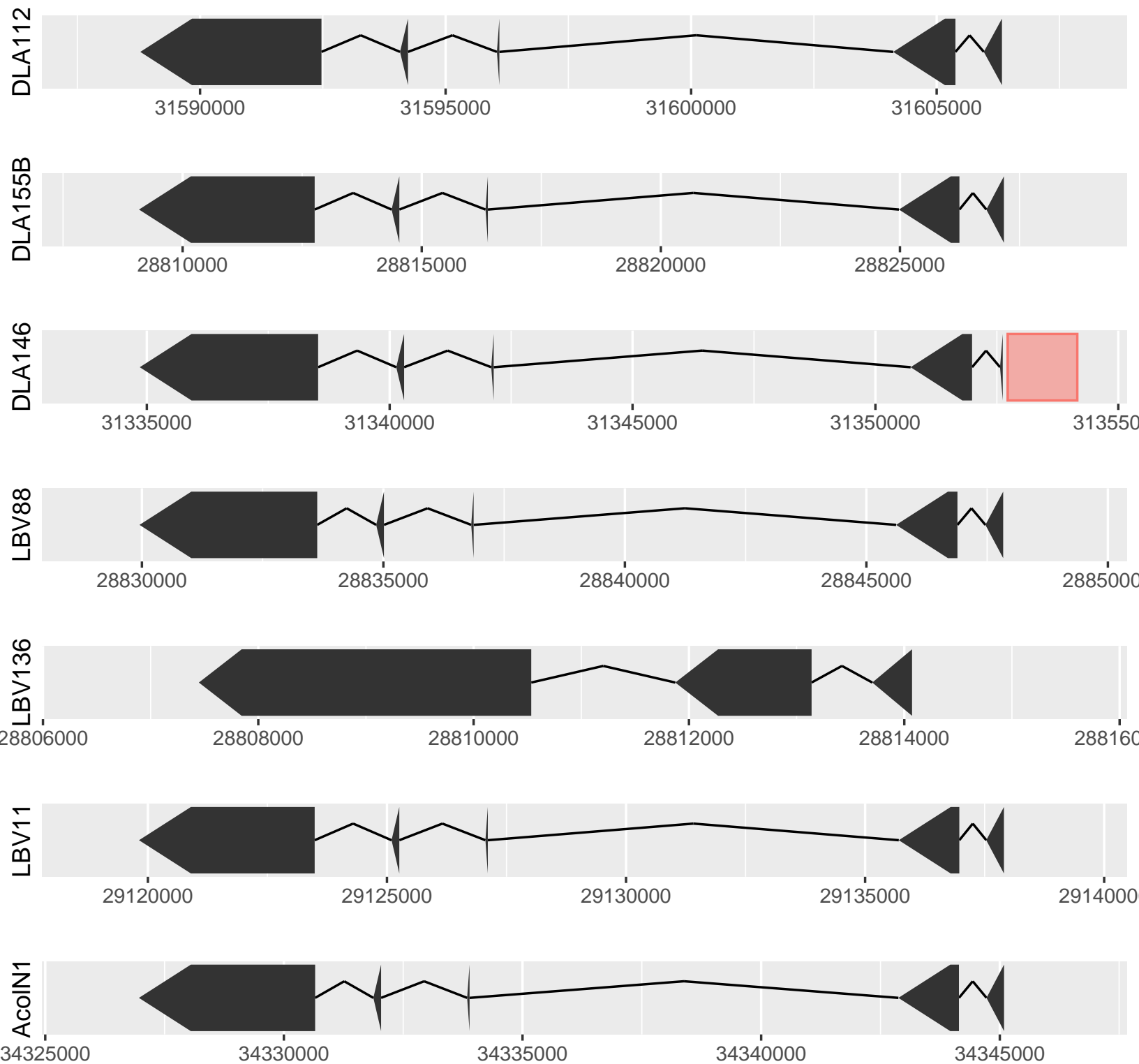
AGAP028069



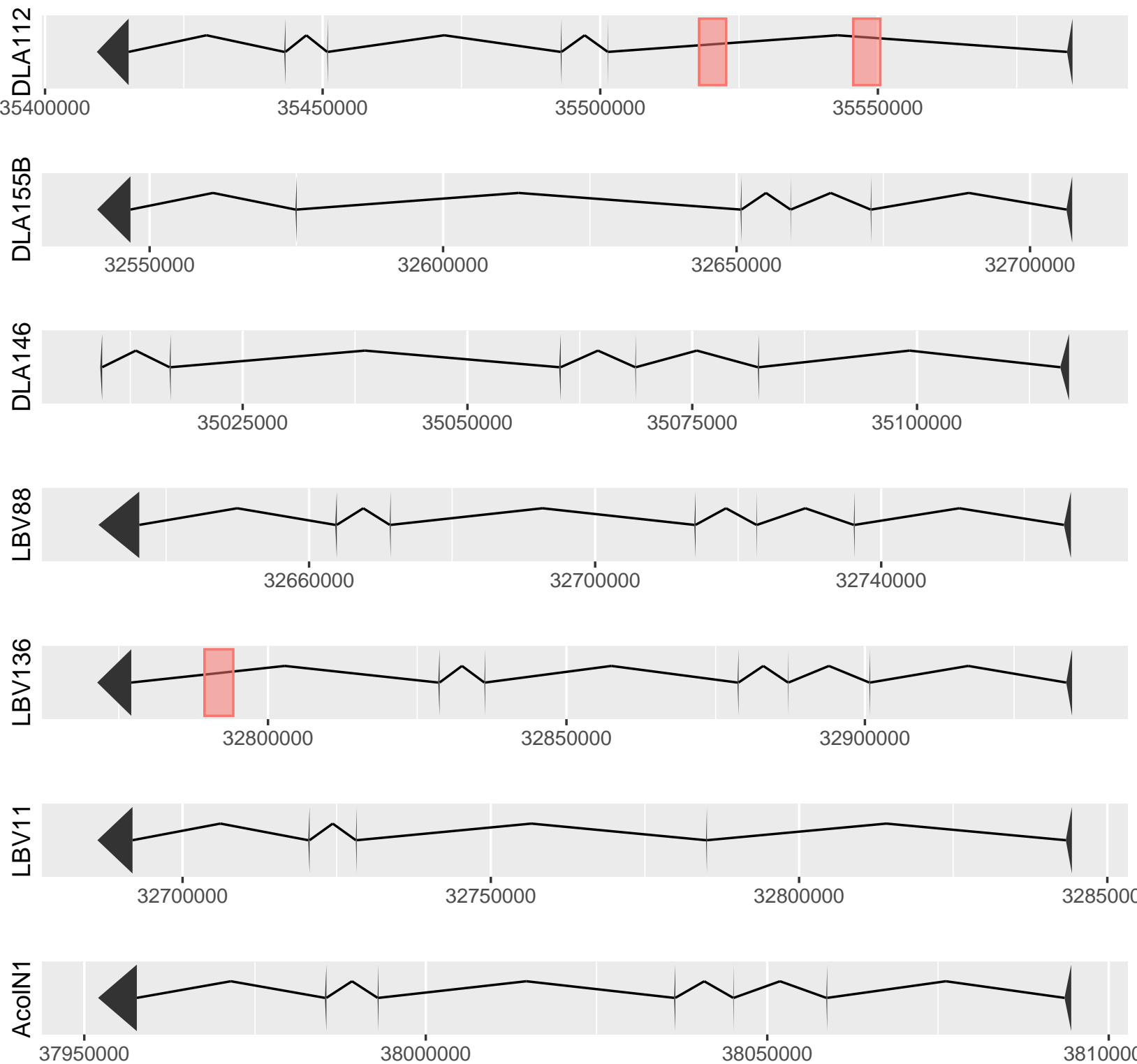
AGAP028449



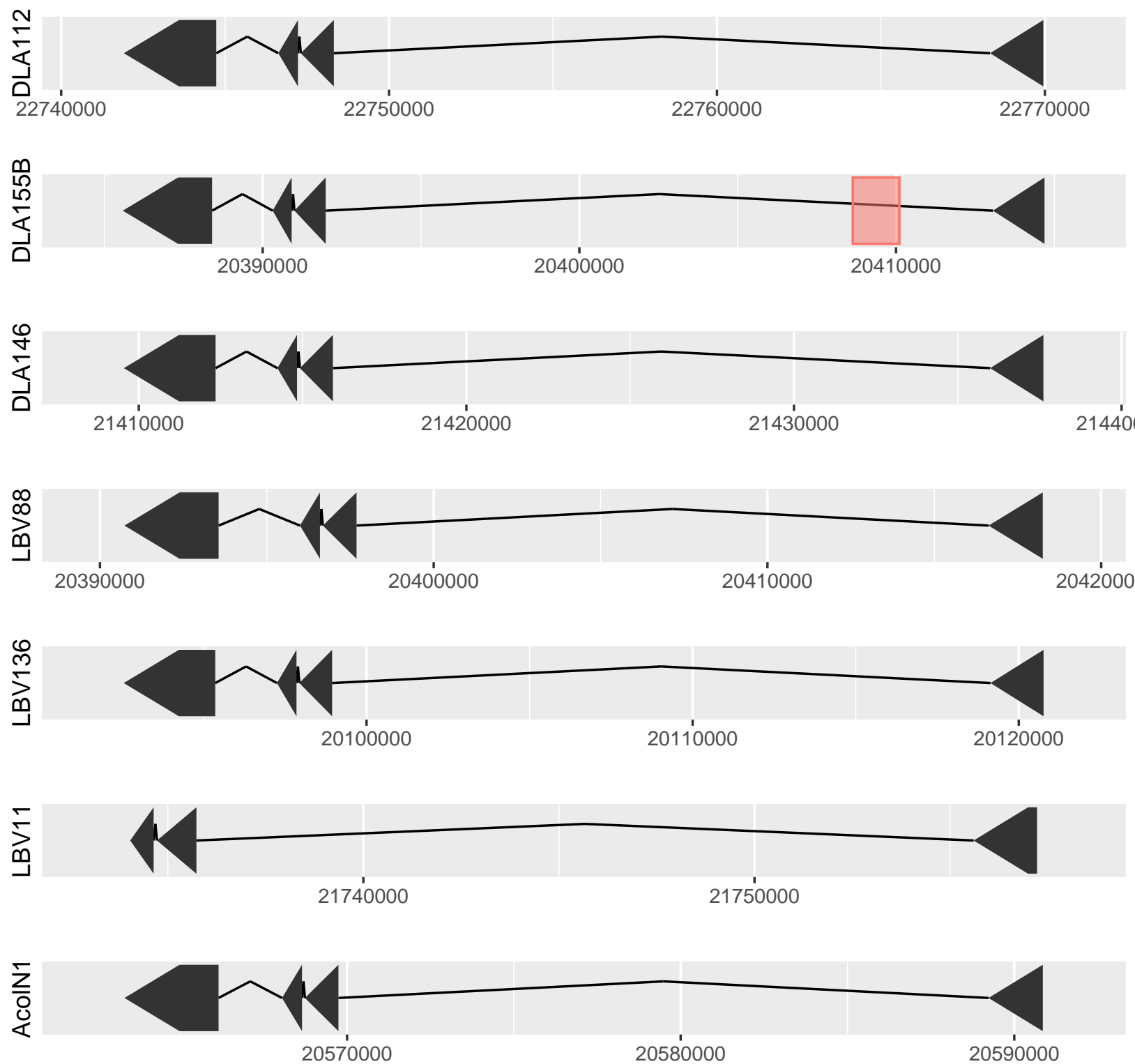
AGAP029185



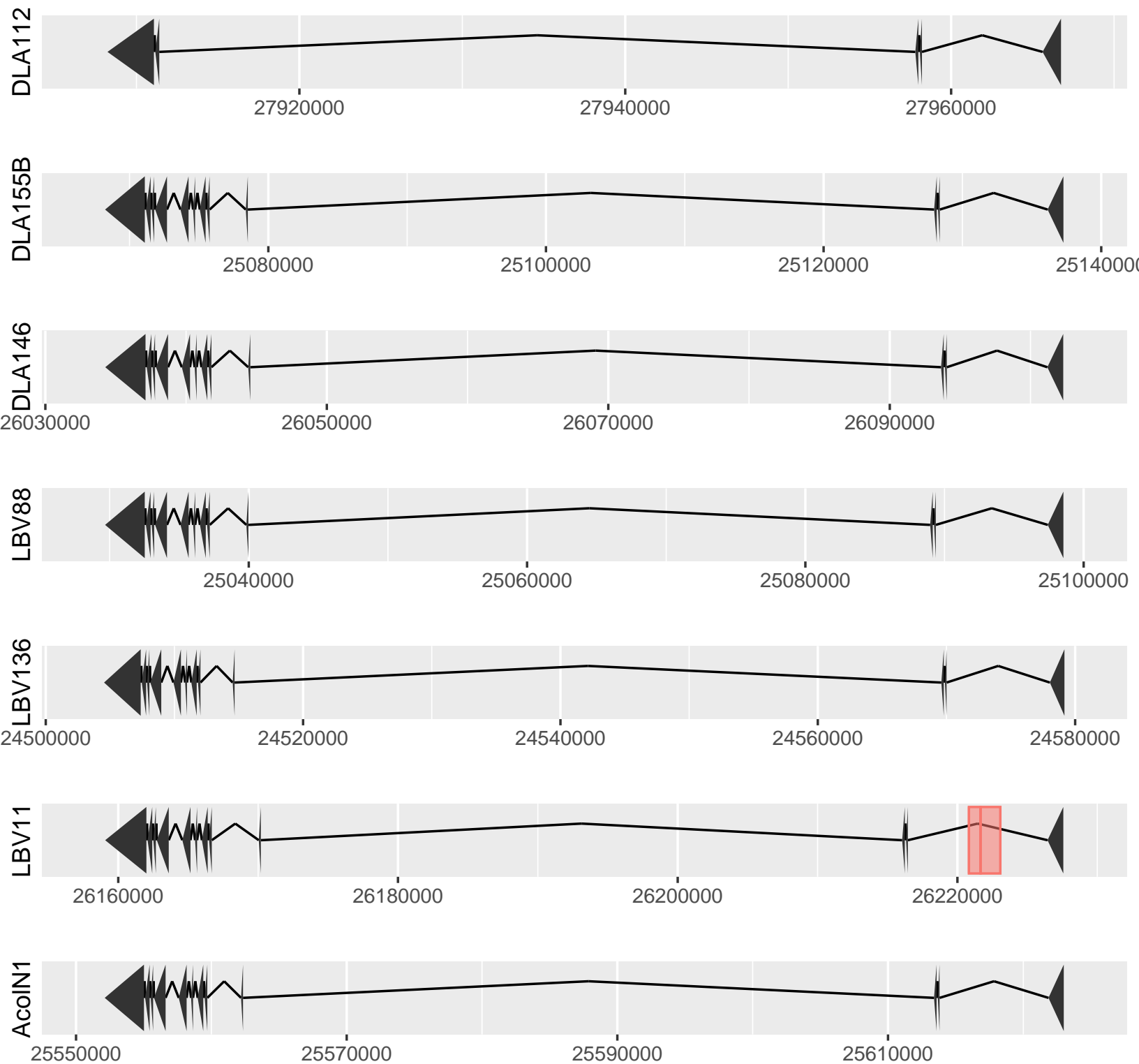
AGAP029191



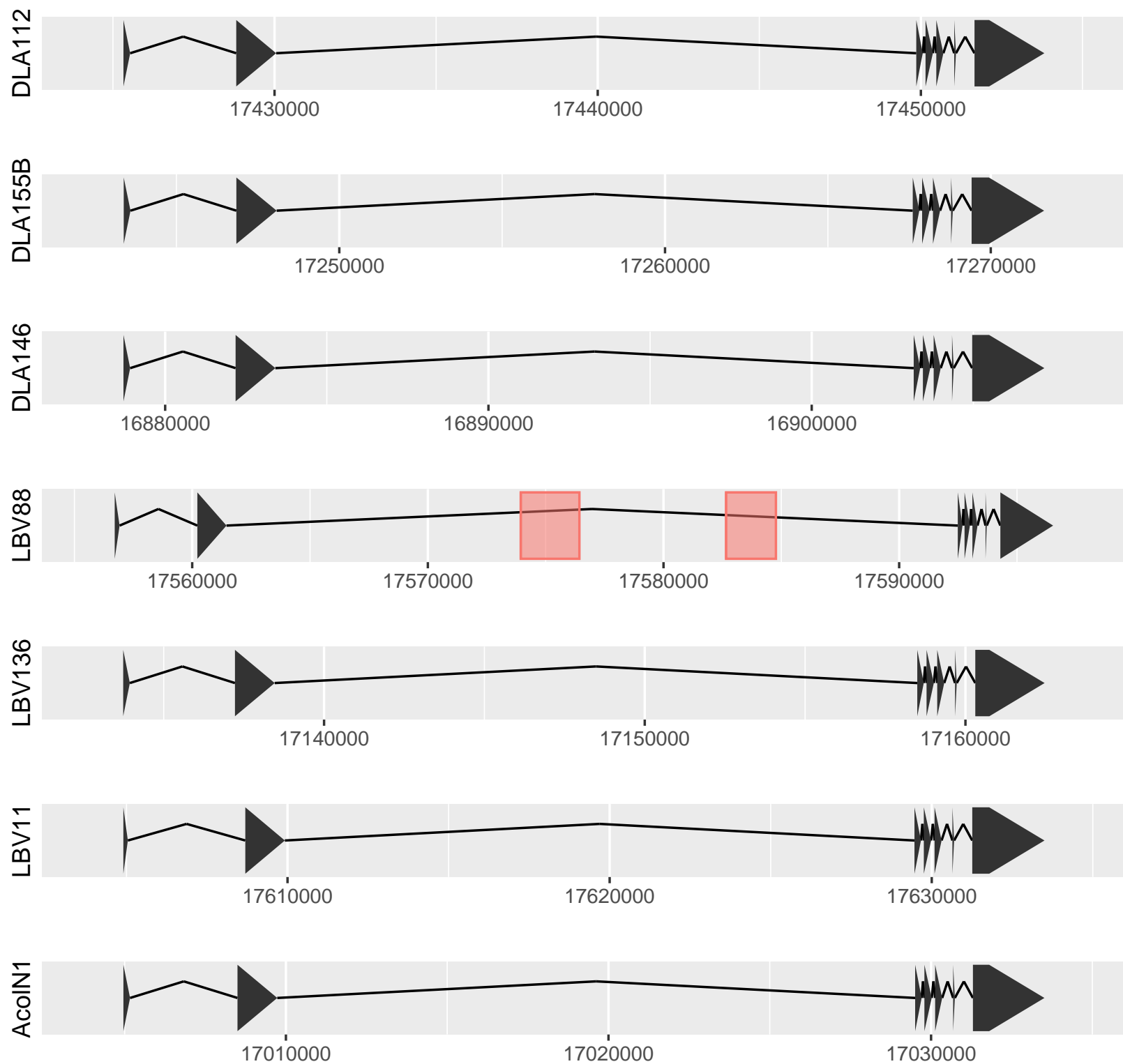
AGAP029293



AGAP029551



AGAP029565



AGAP029620

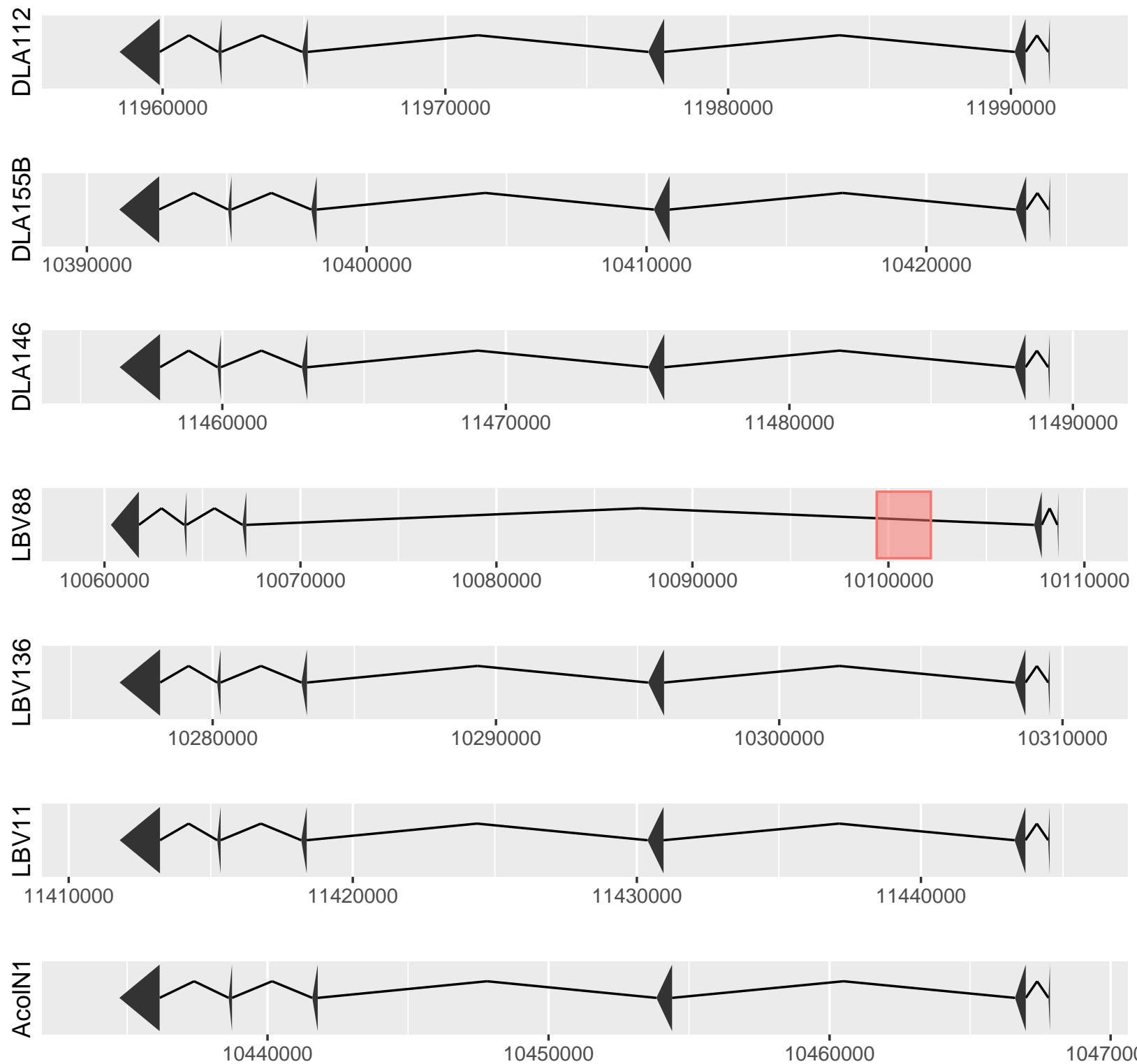
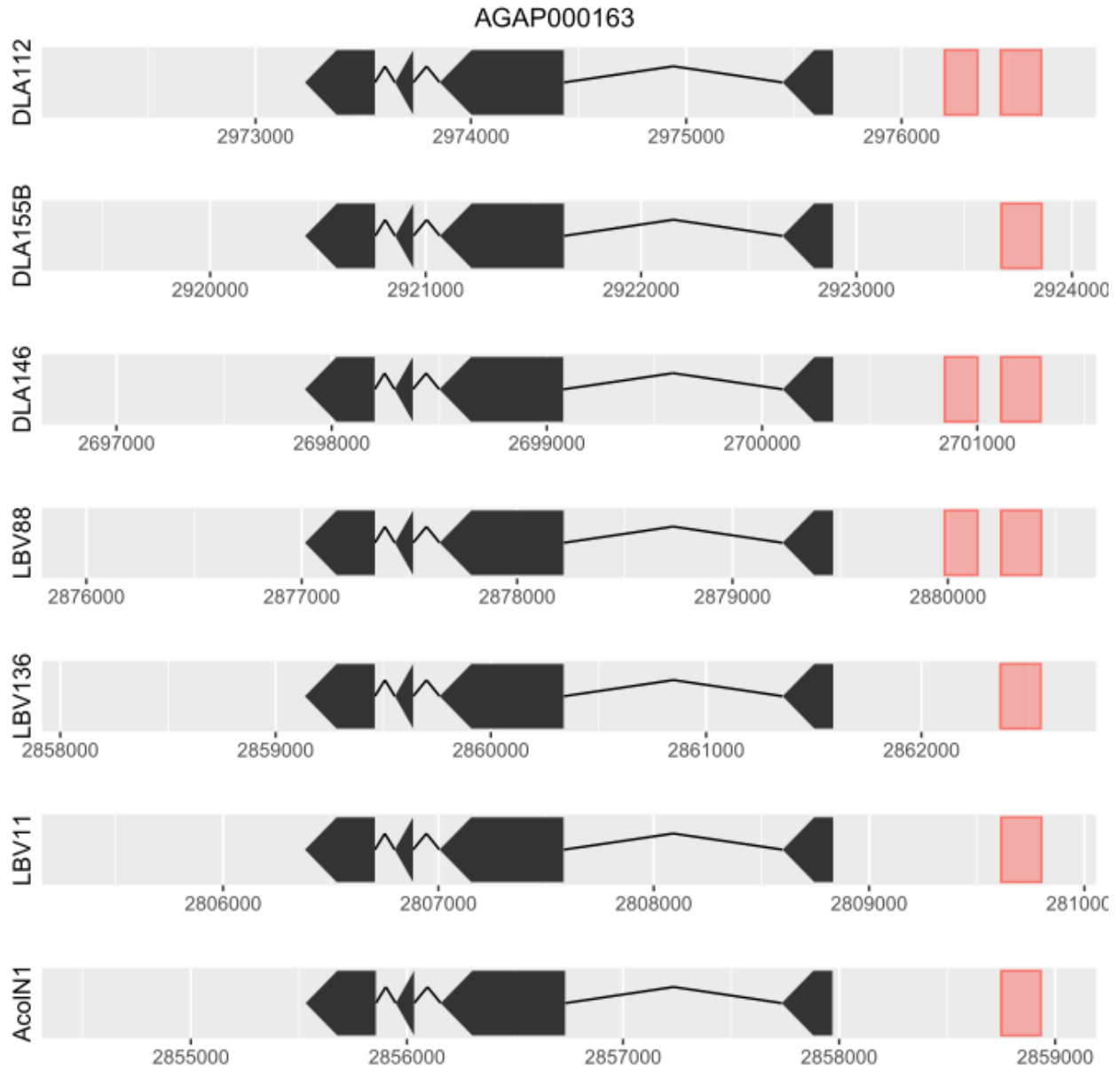
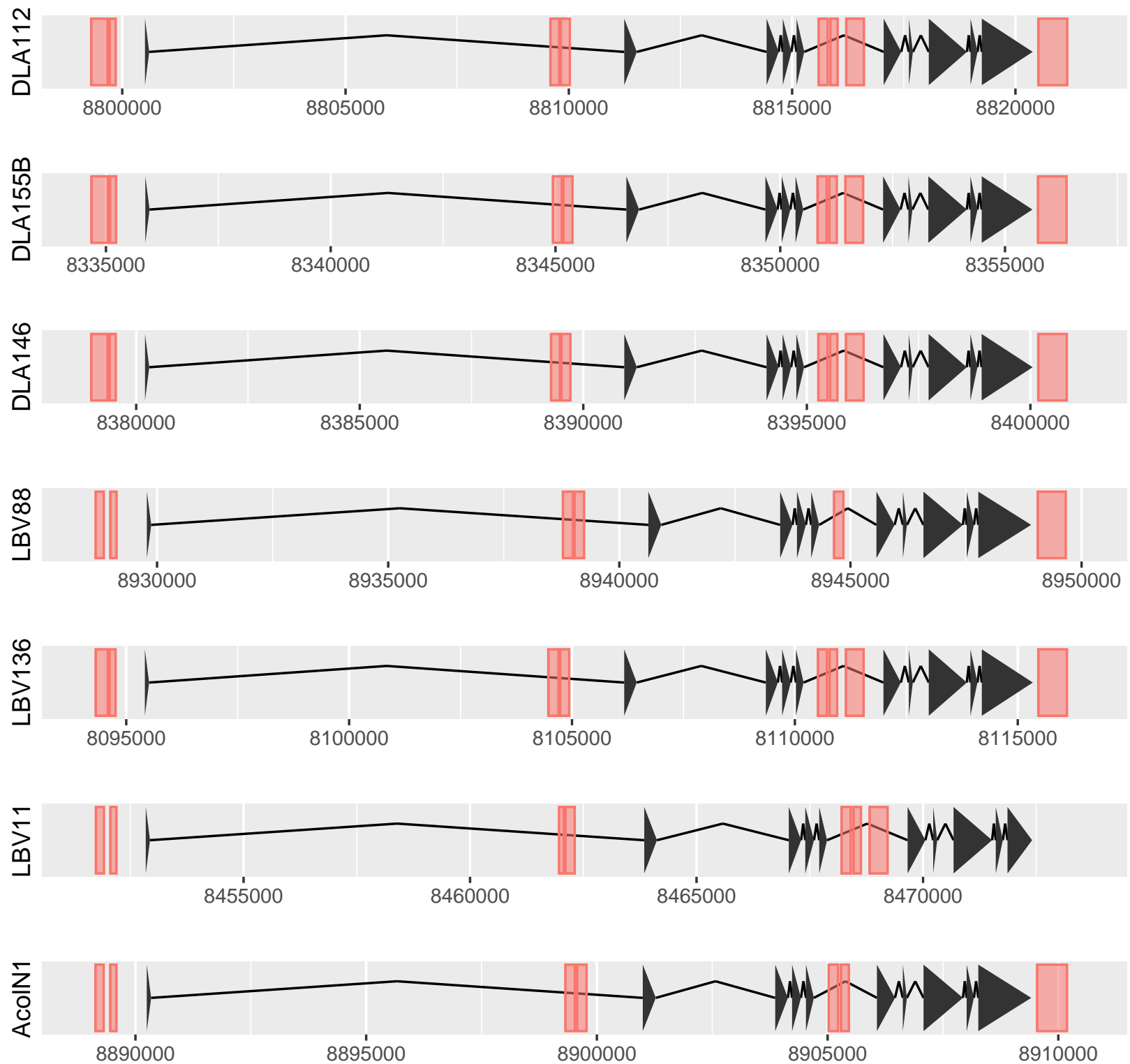


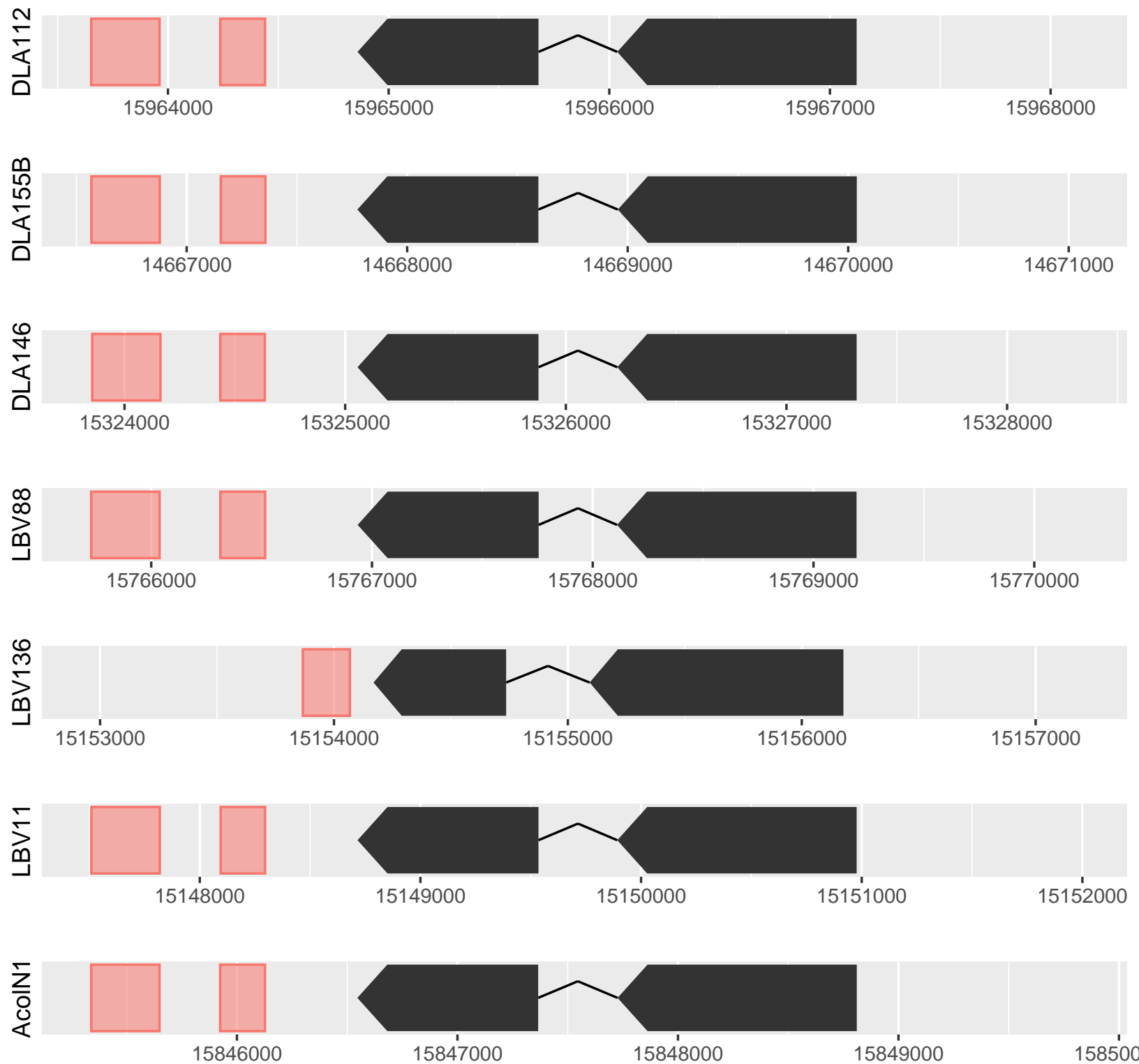
Figure S7. Genes associated with insecticide resistance with TE insertions. Diagrams of genes associated with insecticide resistance showing the gene structure and the TE insertions (in red) inside or closer than 1 kb to gene.



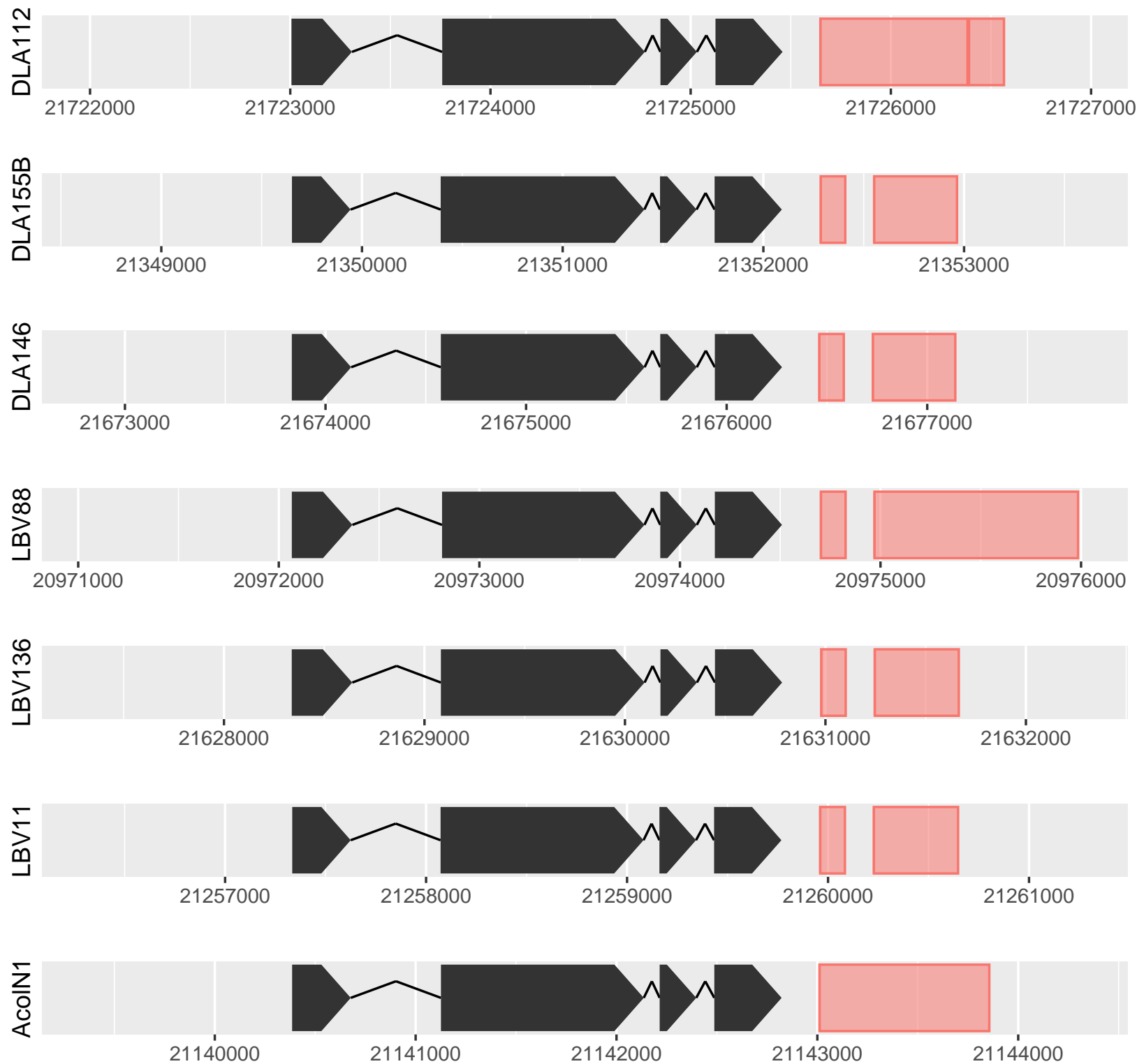
AGAP000500



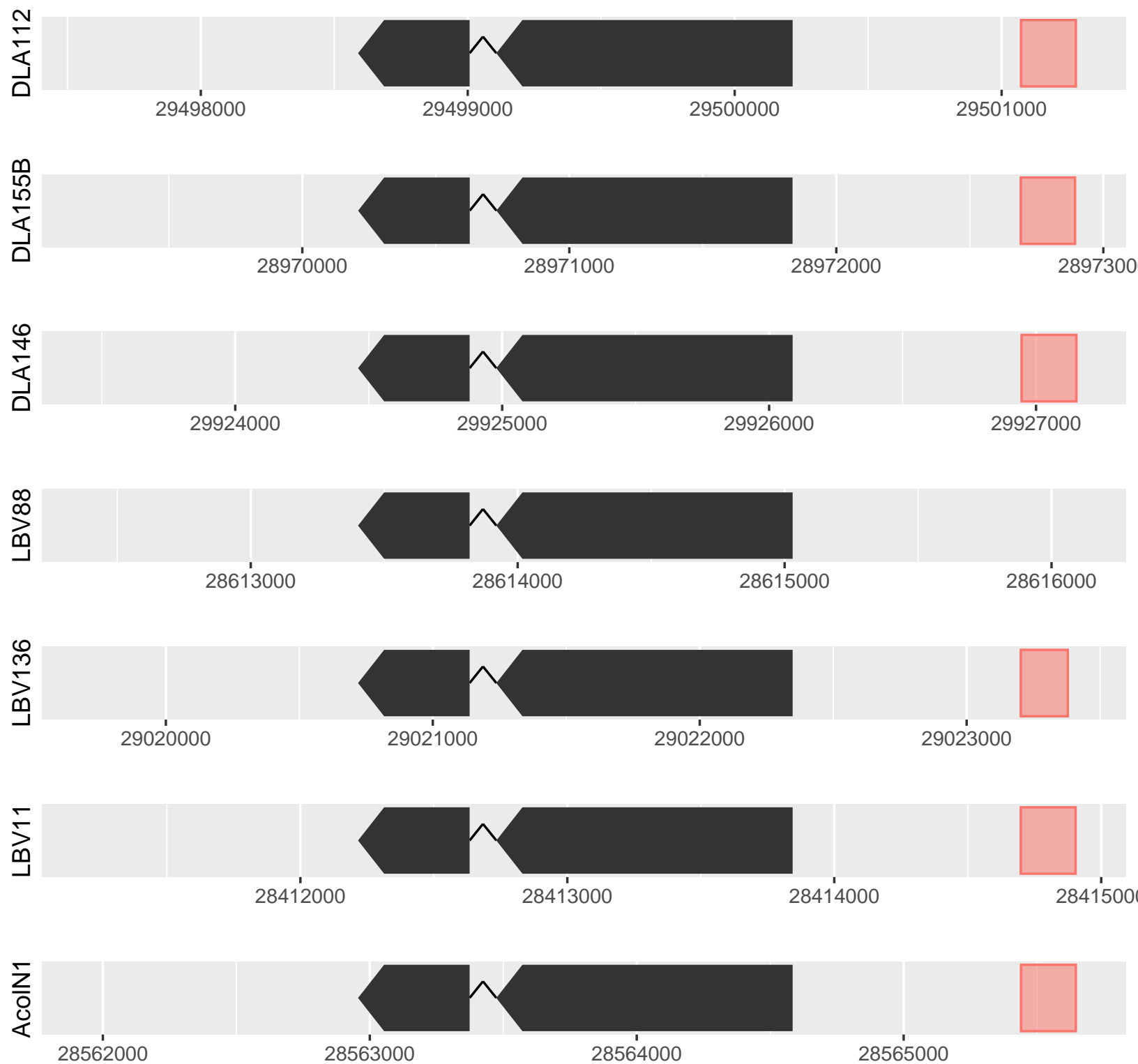
AGAP000818



AGAP002419



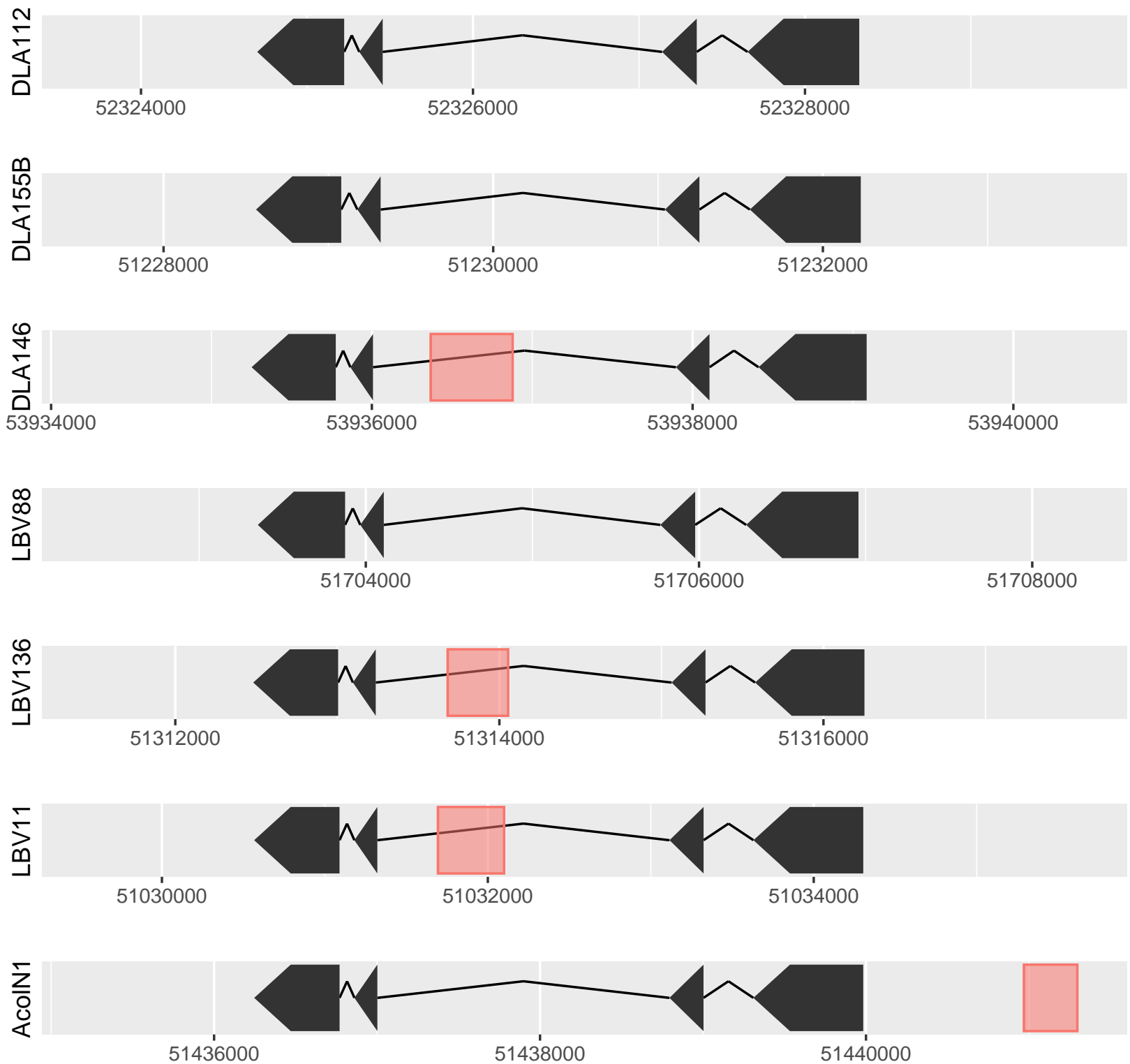
AGAP002866



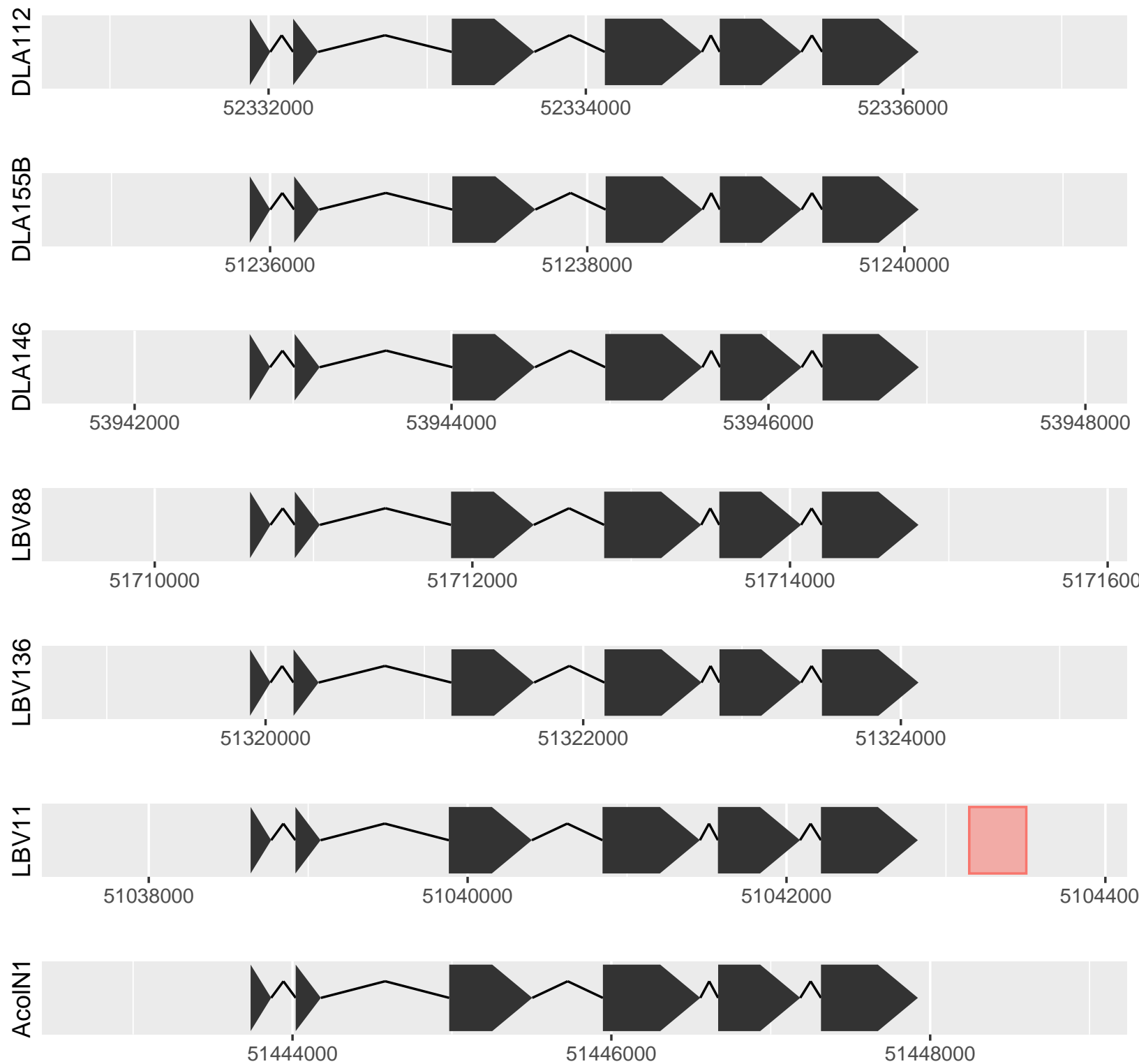
AGAP002867



AGAP004163



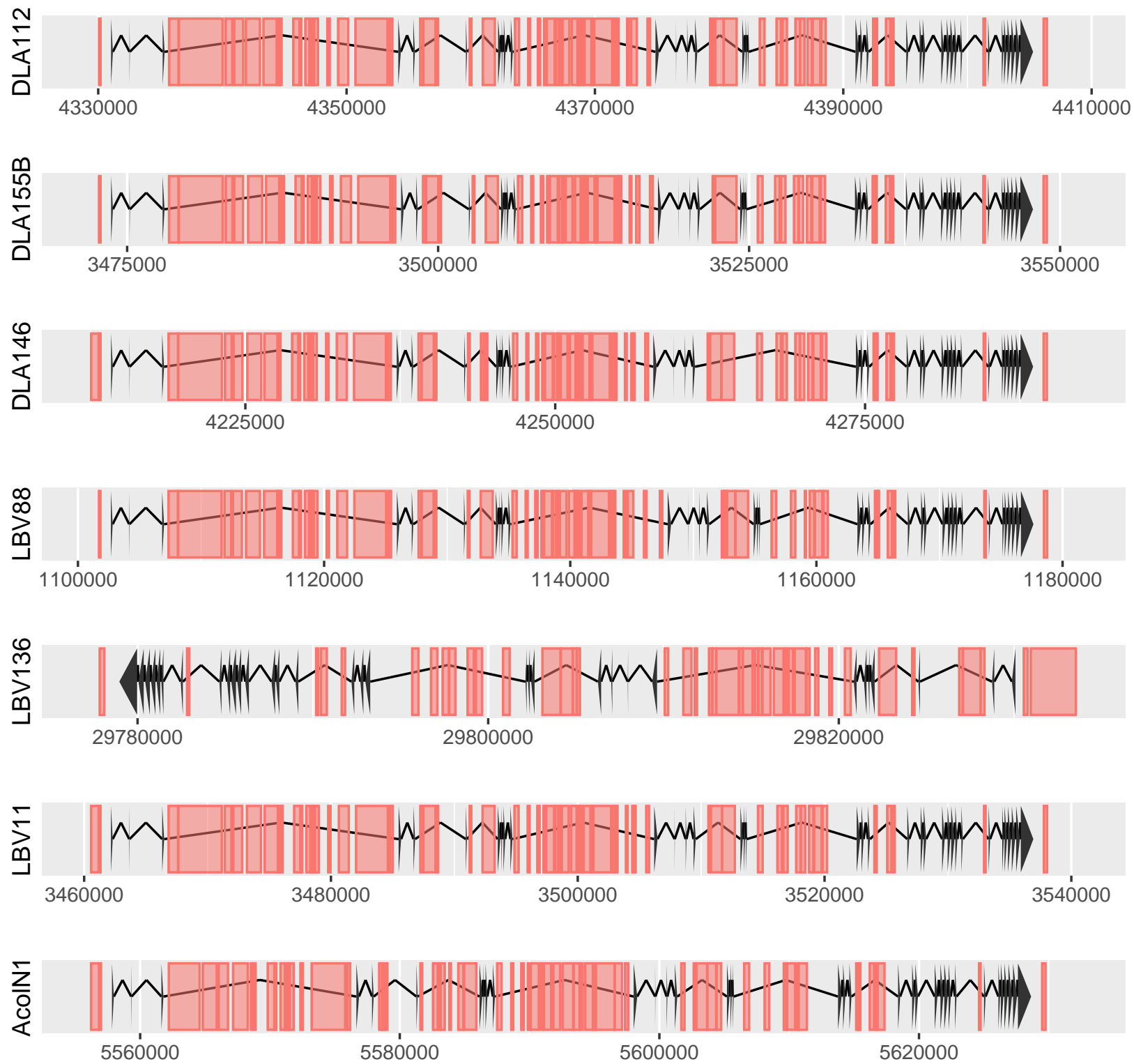
AGAP004164



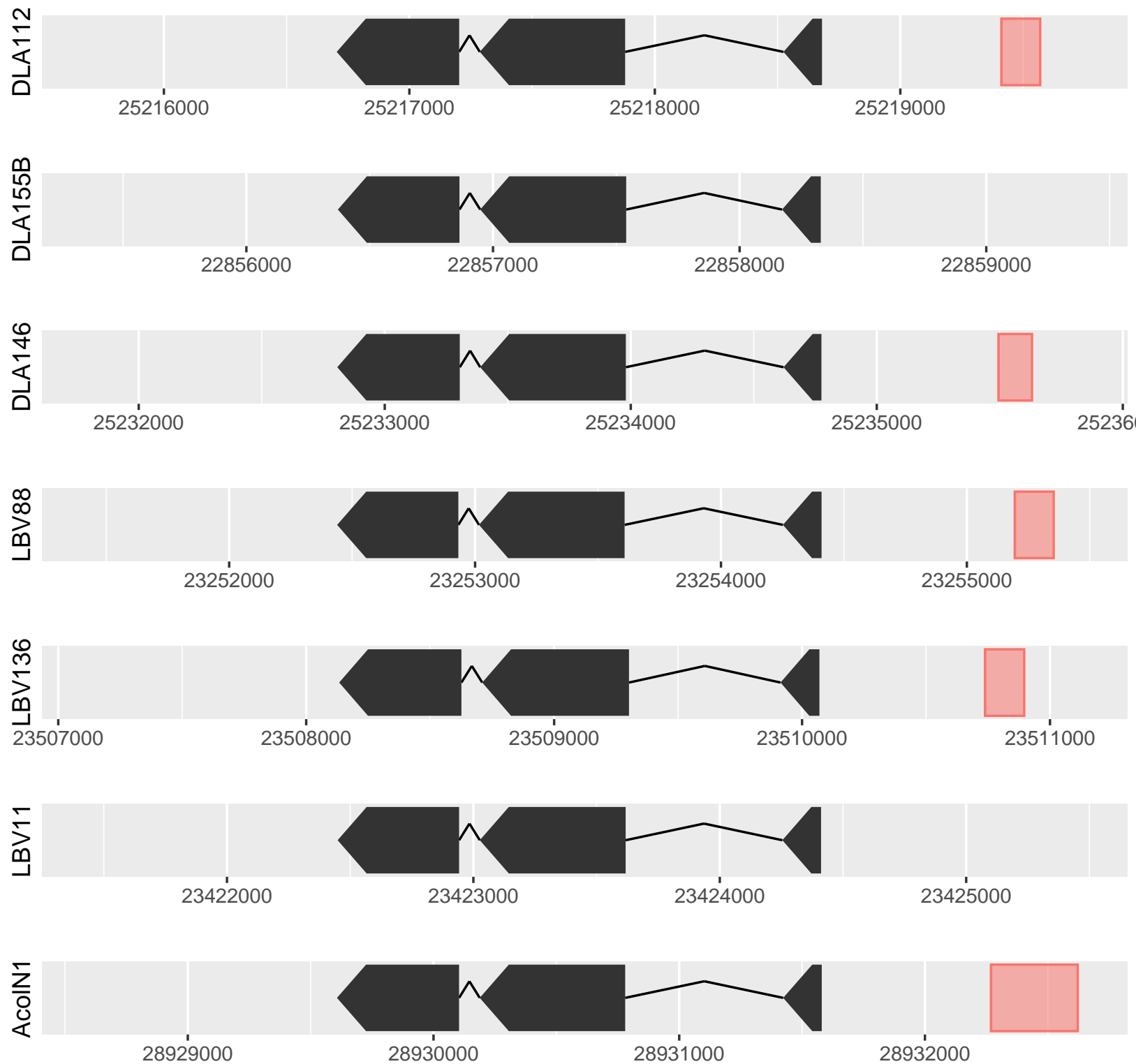
AGAP004382



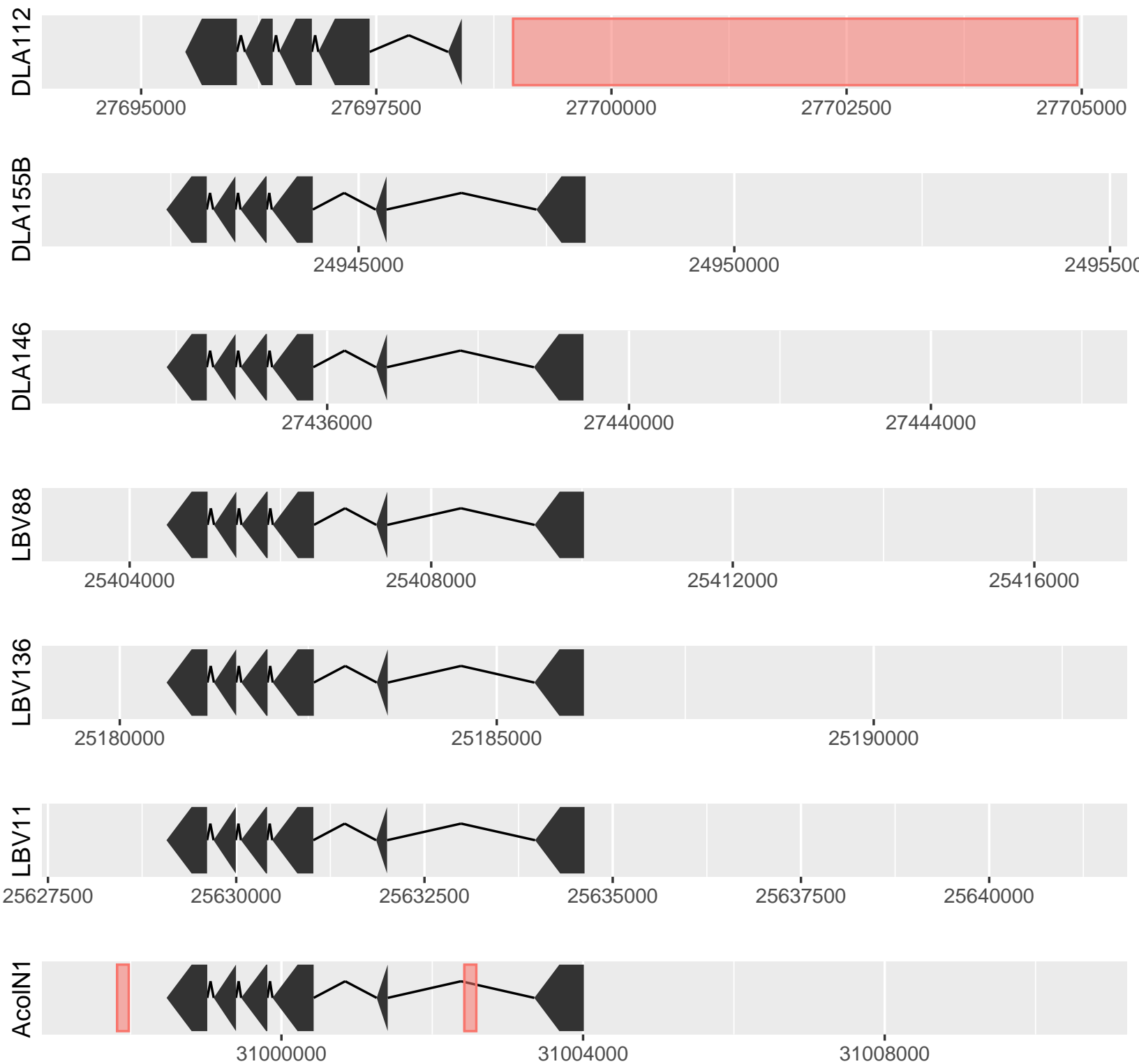
AGAP004707



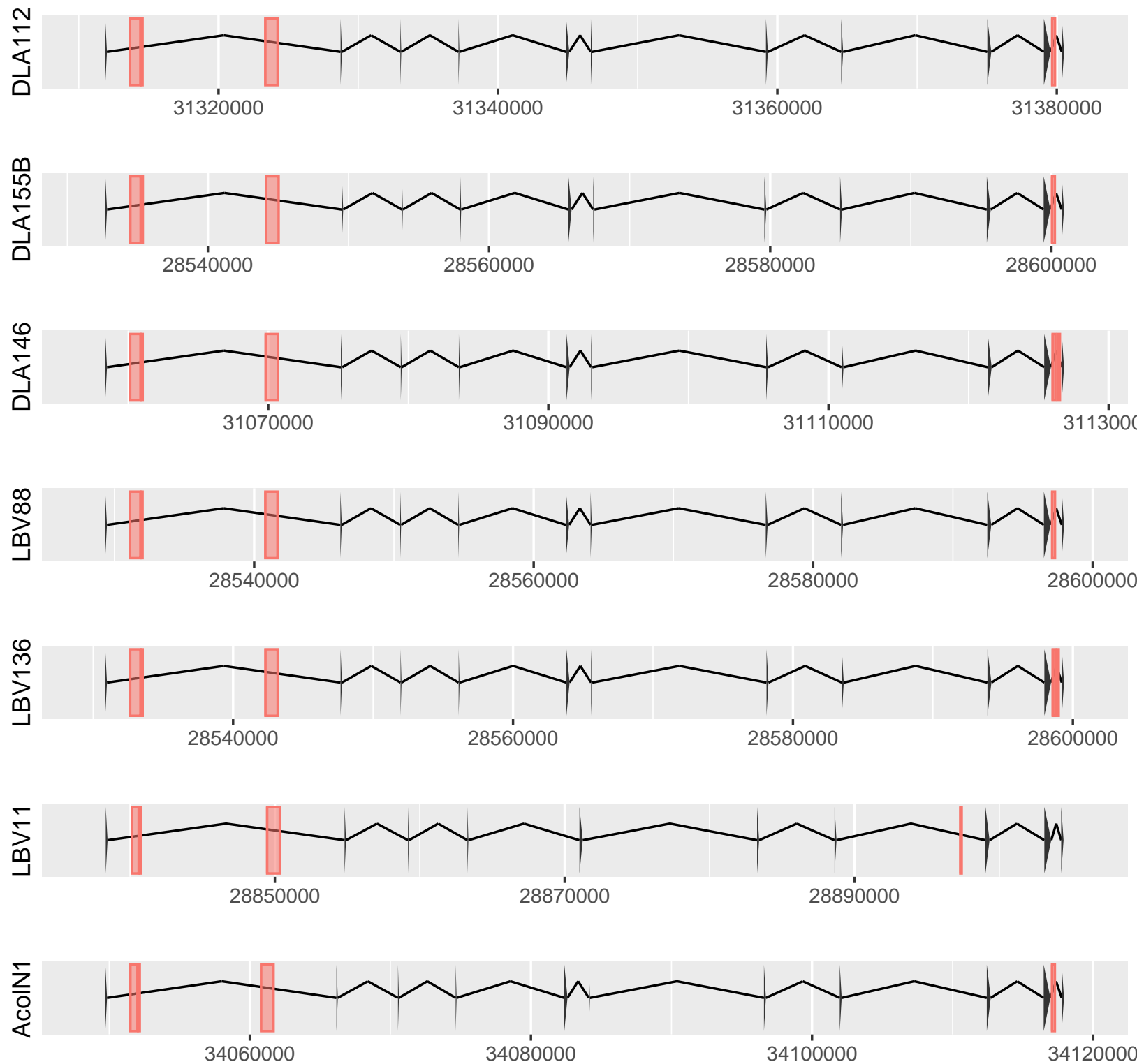
AGAP005749



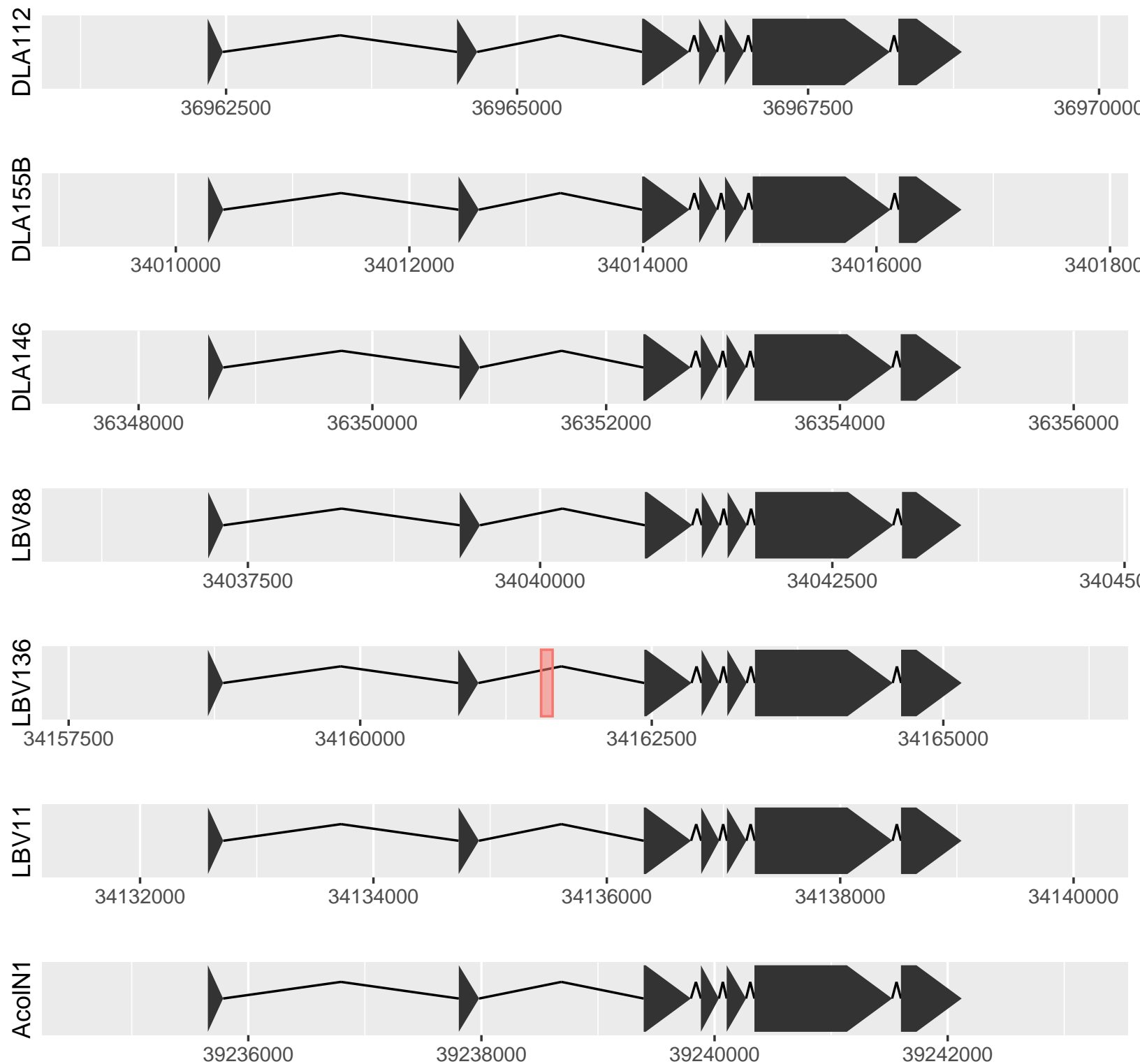
AGAP005834



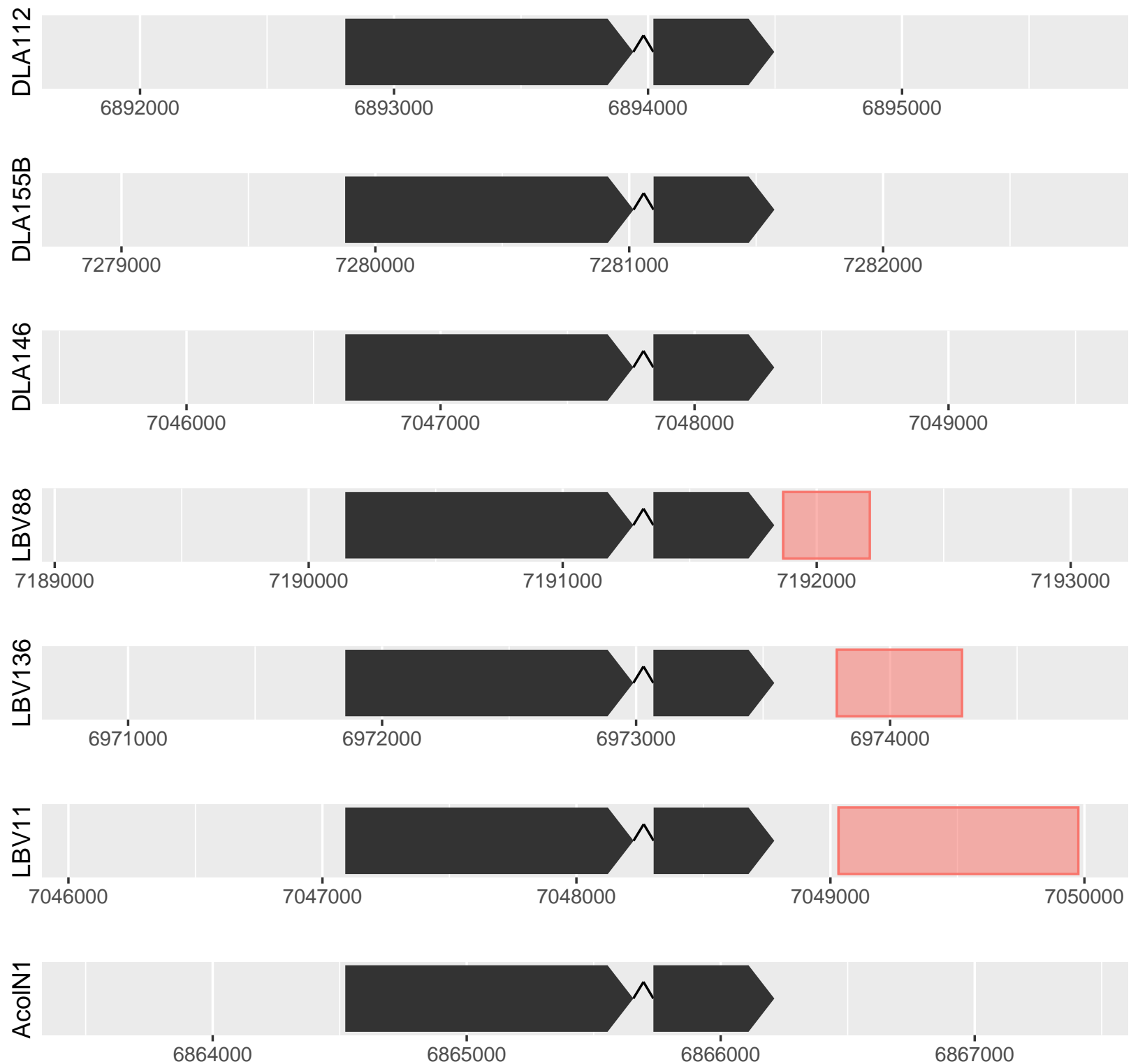
AGAP006028



AGAP006364



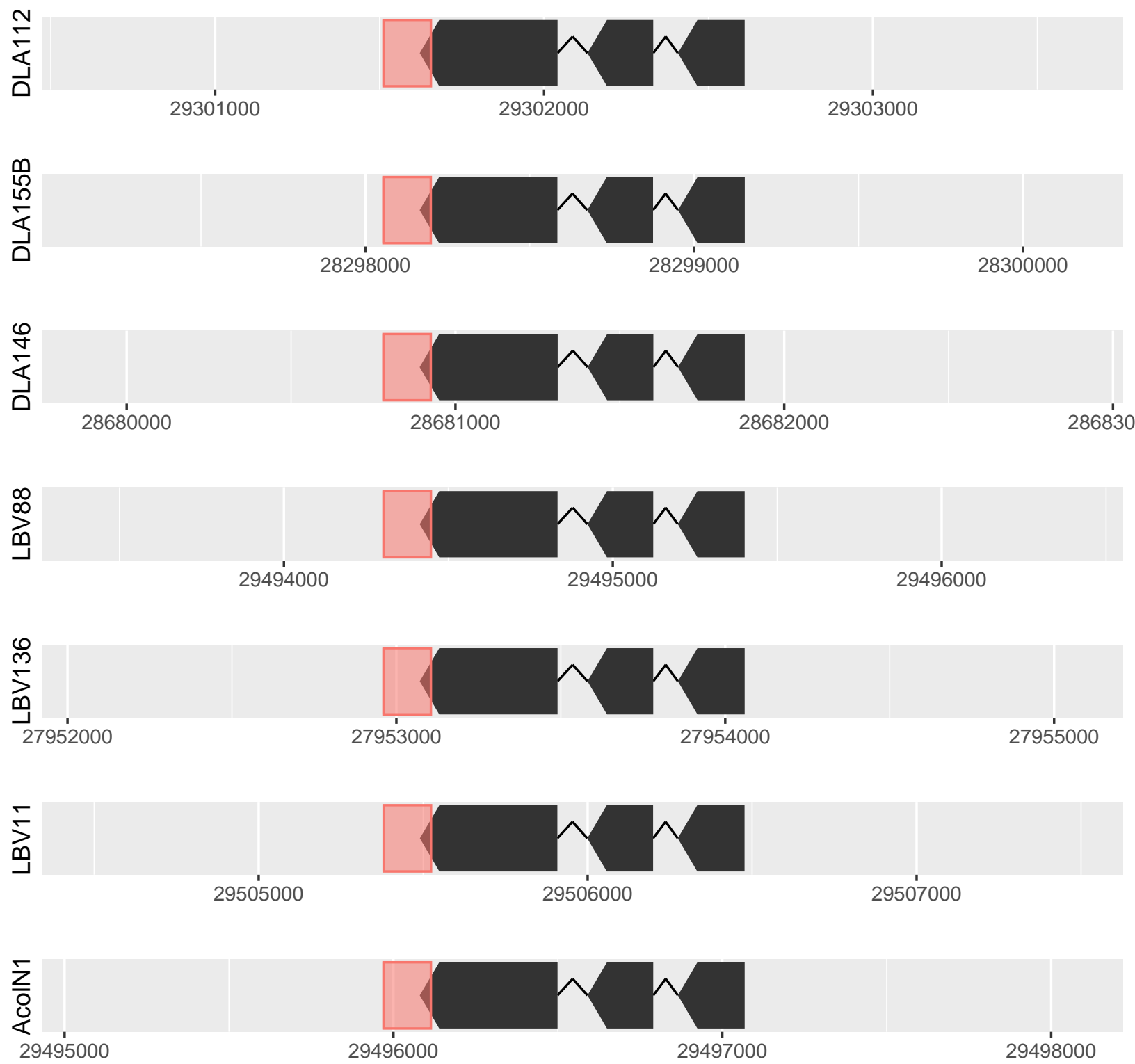
AGAP008212



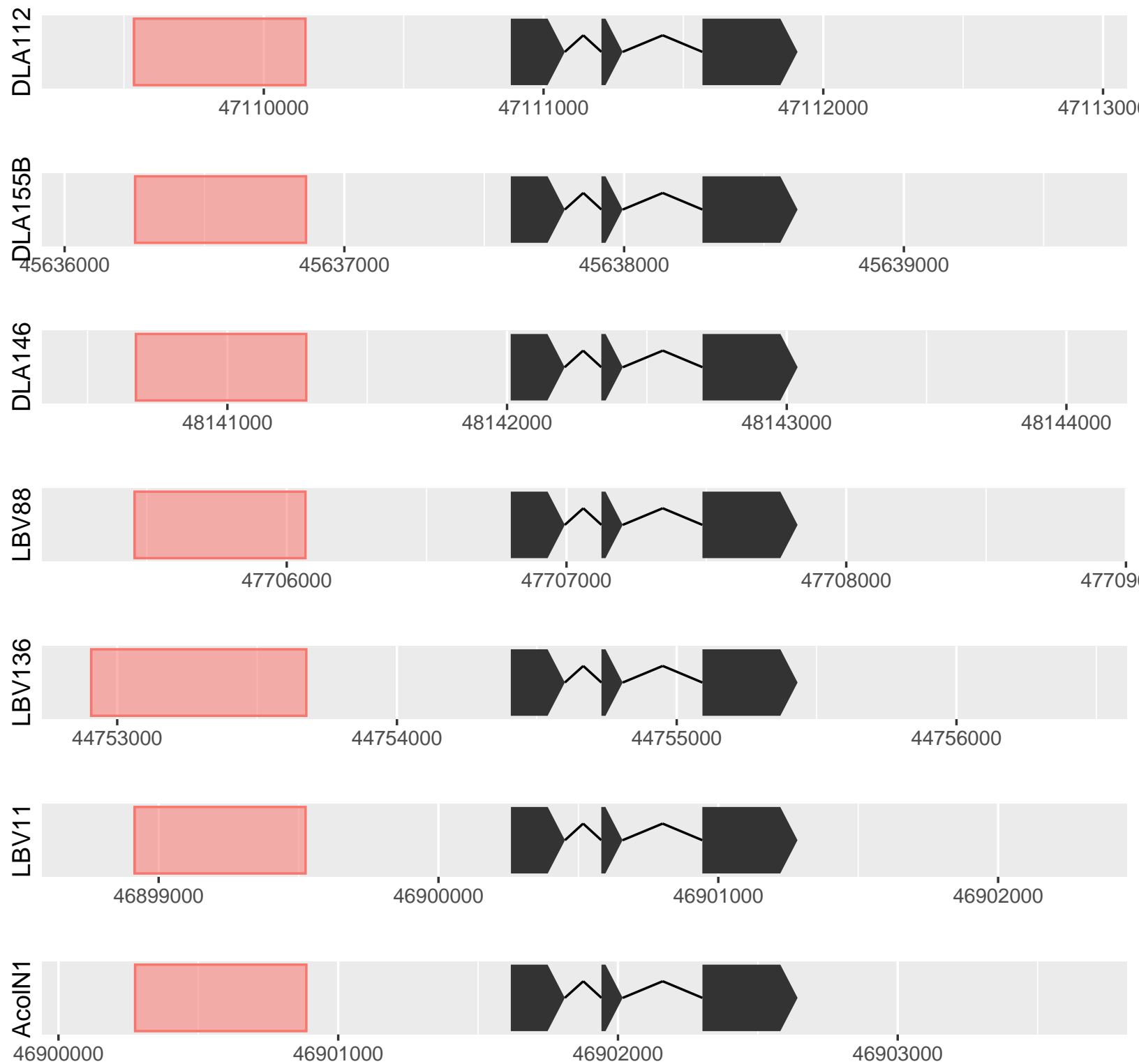
AGAP008213



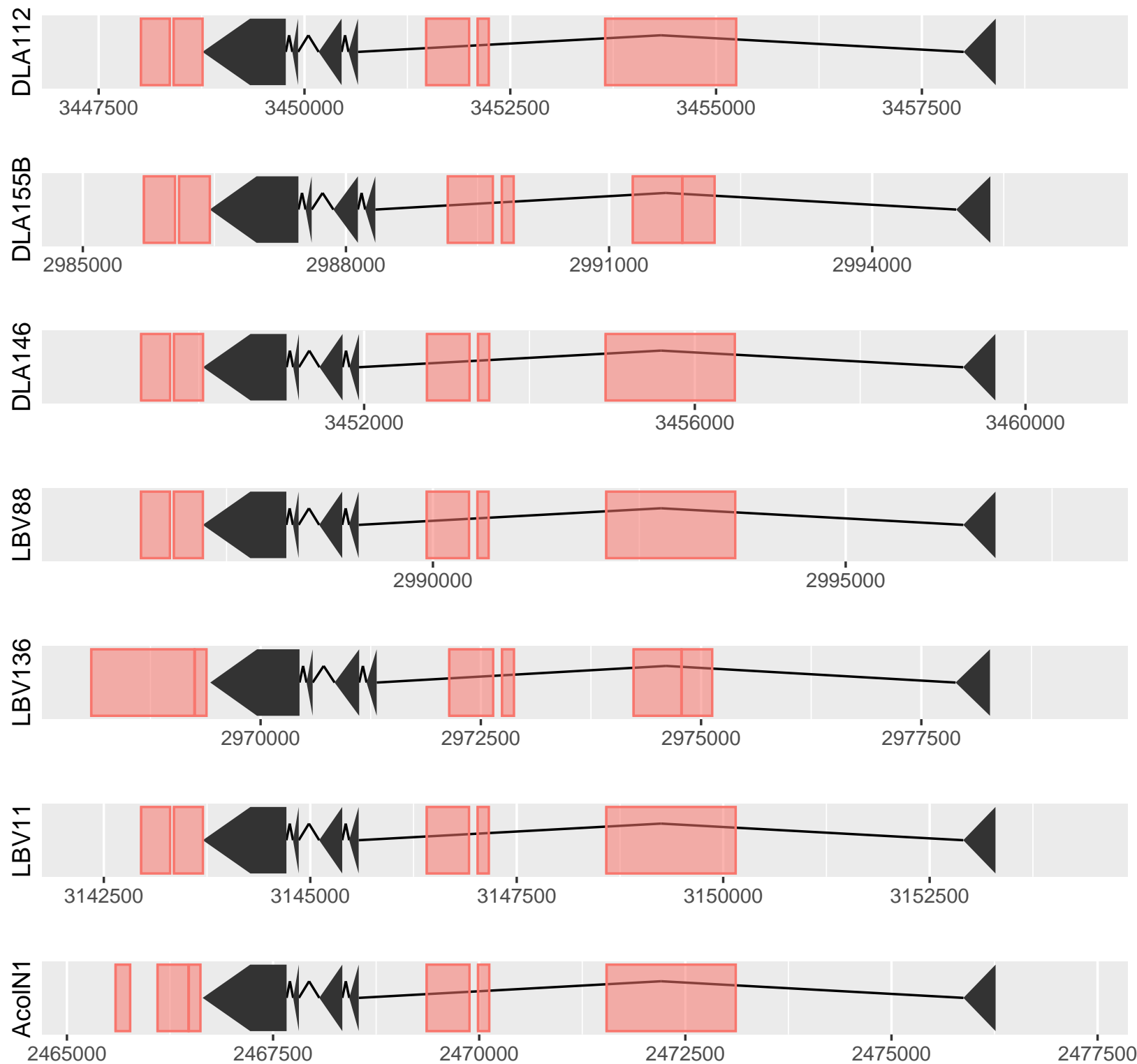
AGAP009194



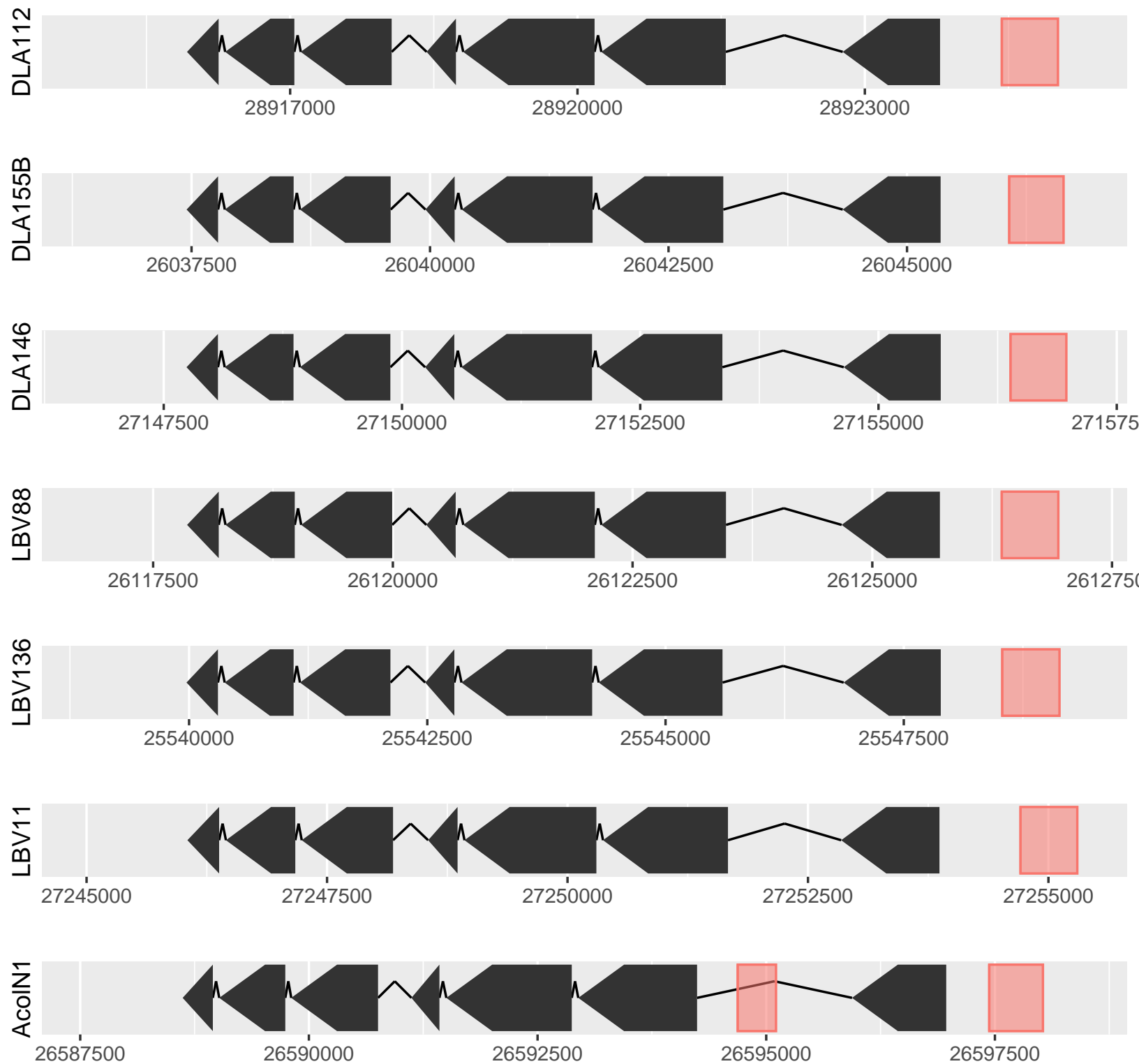
AGAP009946



AGAP010414



AGAP011518



AGAP013121

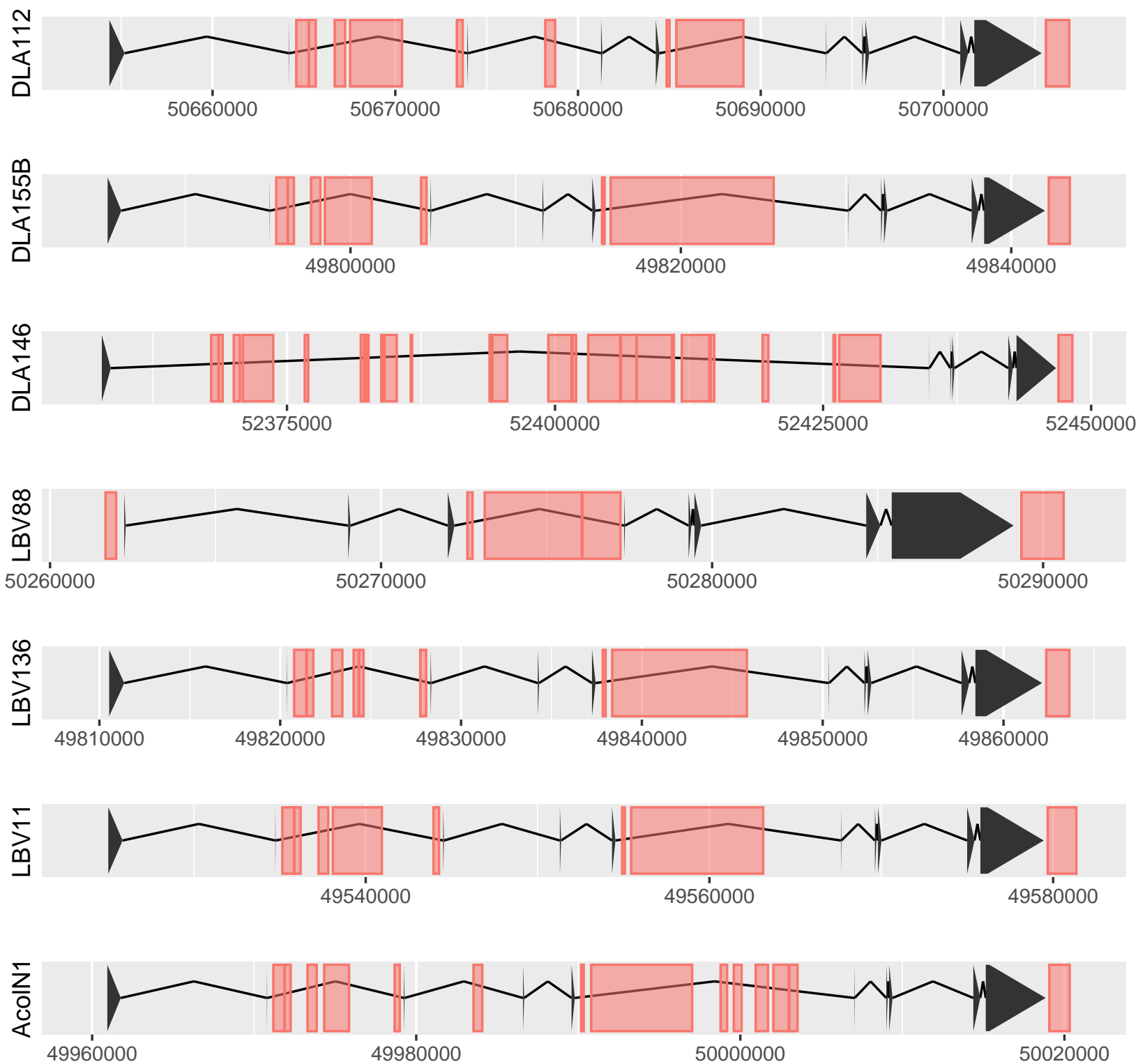
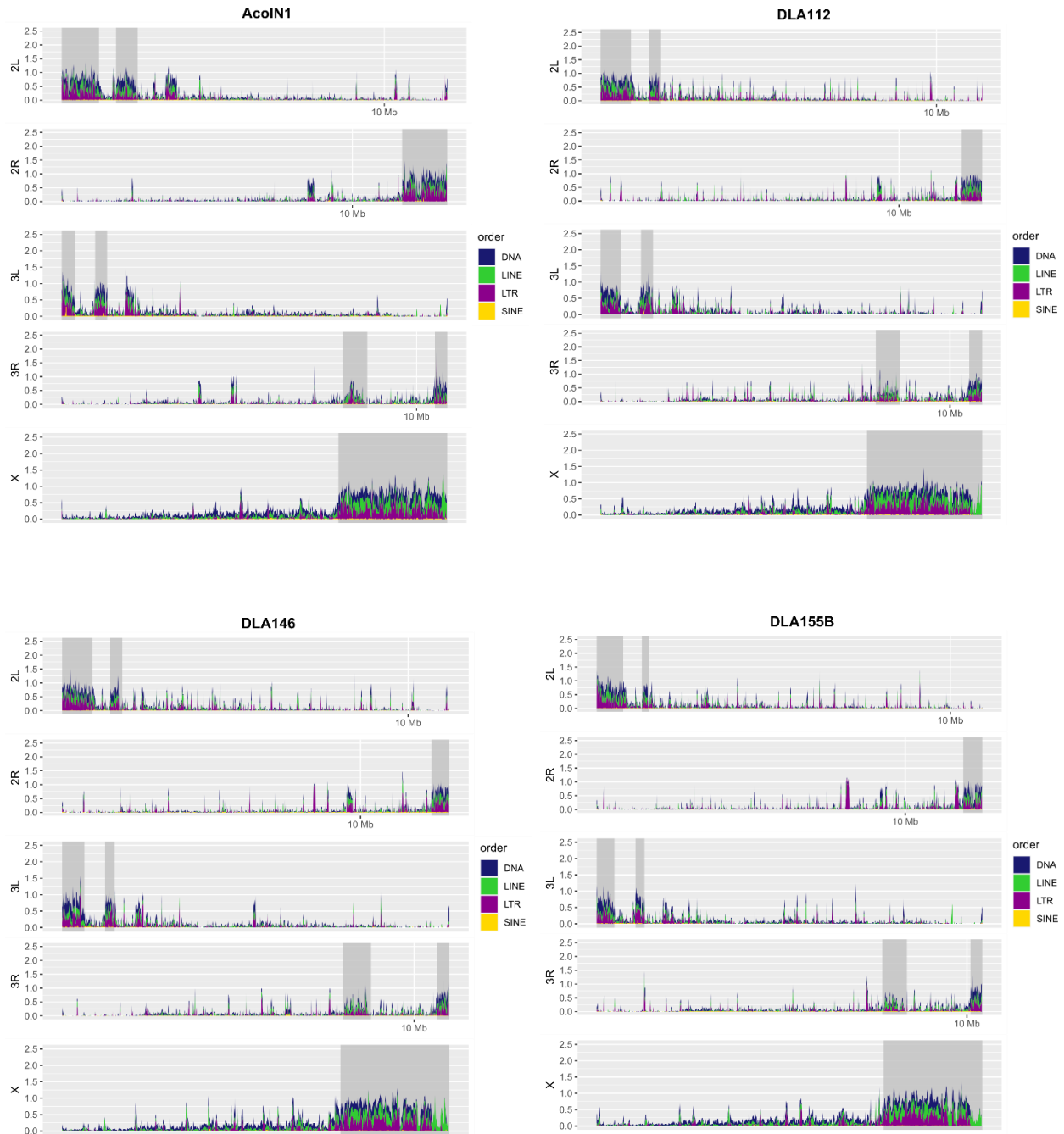
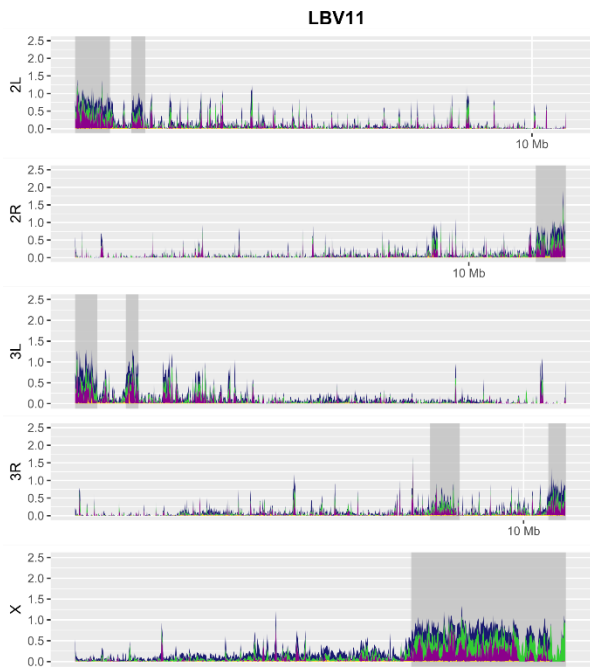
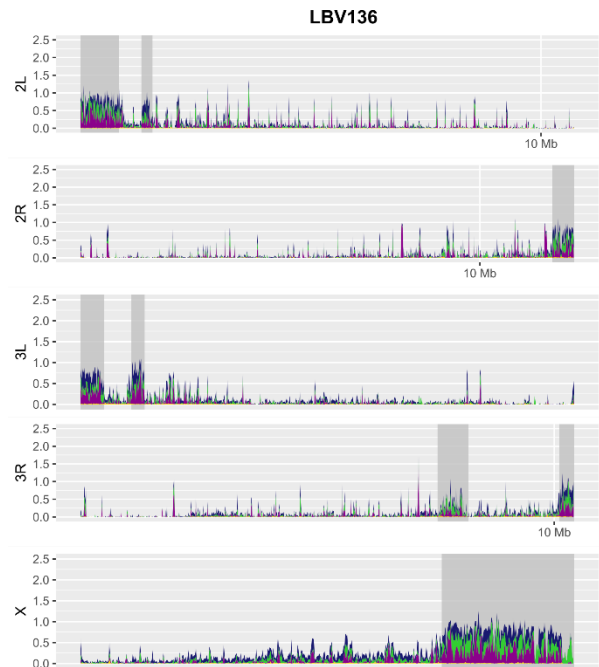


Figure S8. TE abundance plots by chromosome. The TE coverage by position is plotted. Grey boxes delimit the regions where the heterochromatin coordinates were transferred from the AgamP3 genome to each of the genomes used in this work.





Supplemental Table Descriptions

Supplemental Table S1. General statistics, linear models, and euchromatin coordinates for all the genomes analyzed in this work.

A) General statistics for the assembled genomes from the six sampled breeding sites. Whenever available, these statistics are also shown for the reference genome *AcolNI*. Scaffolding was performed using the AgamP4 chromosome level genome assembly. B) Linear models testing for significant associations between long-reads coverage, long-reads N50, number of contigs, N50 of contigs, complete BUSCO and number of transferred genes with different sequencing, assembly, and scaffolding statistics as well as different TE content and TE diversity values. FDR values smaller than 0.05 are highlighted in red. C) Coordinates for the euchromatin transferred using the reference from Sharakhova *et al.*, (2010).

Supplemental Table S2. TEs identified using REPET.

A) Number of TEs identified in each stage of the annotation process. TEs with MAP > 1 correspond to TEs annotated using more than one sequence. TEs with FC >3 correspond to TEs with more than three complete copies in the genome in a single or multiple blocks. TEs with FF >1 correspond to TEs with more than one complete copy in the genome in a single block. Reliable TEs correspond to TEs that upon observation of the coverage plots and the presence of features typical of their orders were considered to be accurate. TE consensus manually corrected correspond to TEs that although reliable required removal of features or extension, this category comprised mostly chimeric annotations. B) Families identified in each genome.

Supplemental Table S3. TE families annotated in *An. coluzzii*.

A) "Known" describes whether the family had previously been described in another species (Yes) or first described in the present study (New). "Classification" describes whether families are degenerate (DG), putatively autonomous (PA), putatively autonomous but lacking TIRs or LTRs (PA*) or putatively non autonomous (PNA). "BBH in TEfam" refers to the ID of the bidirectional best hit in TEfam. This applies as well for BBH in AnnoTEExcel and BBH in Repbase. B) Summary of the number of insertions or bases from each order or superfamily in each of the seven genomes.

Supplemental Table S4. Newly identified TE families.

The 64 families were classified as follows: 32 DNA, 9 LINEs and 23 LTRs. Moreover, the majority of the new families (43/64) had partial matches to other known TEs, thus allowing us to classify them at the superfamily level. For each genome, the total number of insertions, the number of NFL (near full-length fragments, covering at least 90% of the consensus) and FL (full-length fragments, with differences in size with the consensus of less than 30 bp) insertions are given. The median number of insertions considering all strains and the fraction of insertions that are NFL are also given. The number of insertions with greater than 80% identity and covering at least 80% of the consensus in the reference genomes of several species of the *Anopheles* genus are shown and their columns colored differently according to the series they belong to. Aalb - *An. albimanus* (Nyssorhynchus), Aatr - *An. atroparvus* (*Anopheles*), Afar - *An. farauti* (Cellia [Neomyzomyia]), Afun - *An. funestus* (Cellia [Myzomyia]), Aste - *An. stephensi* (Cellia [Neocellia]), Aepi - *An. epiroticus*, Achr - *An. christyi*, Amer - *An. merus*, Agam - *An. gambiae*, Acol - *An. coluzzii*, Amel - *An. melas*, Aara - *An. arabiensis*, Aqua - *An. quadriannulatus*, Abwa - *An. bwambae* and

Afon *An. fontenillei* (Cellia [Pyrethophorus]), Cqui - *Cx. quinquefasciatus*, Aaeg - *Ae. aegypti*, Dmel - *D. melanogaster*. The five anopheline species that do not belong to the Pyrethophorus series were selected given that they had the most complete genomes in their taxonomic group.

Supplemental Table S5. Distribution of TE insertions regarding genes. **A)** χ^2 tests for the total number of insertions over the size of each compartment. The expected values are calculated determining the number of insertions that would be expected given the size of the compartment. The residuals for each category are shown. **B)** χ^2 tests for the number of insertions group by order in each compartment by strain. The residuals are shown, in yellow positions with absolute values larger than 5 and in red positions with absolute values larger than 10, and thus with more weight. **C)** Shows the TEs grouped by superfamily.

Supplemental Table S6. Inversion breakpoints in *An. coluzzii*. The known coordinates or the breakpoints of five inversions were transferred from the *An. gambiae* genome using the Corbett-Detig *et al.*, 2019 reference. Given that the distal 2Ru breakpoint is shared with the 2Rc proximal and 2Rd distal we were able to include it as well but we did not take into consideration the proximal 2Ru given that the exact coordinates for this breakpoint are not known. The reference column includes information of the experimental characterization of the inversion.

Supplemental Table S7. Identification and analysis of active TE families. **A)** Strains by dnaPipeTE refers to the number of strains where more than 50% of the insertions are >99% identical to the consensus according to dnaPipeTE. Identical using BLAST is the median percentage of insertions >99% identical to the consensus in all seven strains. Strains with >2 identical FL is the number of strains with more than two identical full-length insertion. Left and right side LTR are the coordinates of the identified LTRs on the consensus sequence, ORF1 and ORF2 are the identified open reading frames in each family. TSD describes the identified length of the target site duplication. **B)** Identities between both LTRs in full-length insertions of the six putatively active LTRs.

Supplemental Table S8. TE insertions from active families in proximity of genes. **A)** Presence or absence of the active insertions by genome is shown. **B)** Genes with TE insertions from active families. In yellow the six genes selected for further analysis. Status refers to whether the insertion(s) found on each of the genes is (are) fixed, polymorphic or found on a single genome.

Supplemental Table S9. Transcription factor binding sites (TFBS) and promoter prediction for TEs belonging to the active families and insecticide resistance genes. **A)** Position and score for each identified binding site of the transcription factors *cnc:mafS*, dl and *Stat92E*. **B)** Summary of the TFBS identified per genome. **C)** Promoter motifs in the TE insertions in genes involved in insecticide resistance and in insertions from active families located in gene upstream regions. In orange insertions that contained one of the promoter motifs (TATA box, Drosophila Initiator or Drosophila TCT), in blue insertions that contained a pair of the promoter motifs (TATA and either Initiator or TCT) and the distance between these two motifs was the expected one (Sloutskin *et al* 2015).

Supplemental Table S10. TEs in genes which are differentially expressed upon exposure to insecticides. In yellow eight genes well-known to be involved in insecticide response.

Supplemental Table S11. TE frequency estimation in rural populations using PoPoolationTE2. **A)** Accession number of the 37 genomes from Burkina Faso and the 30 genomes from Cote d'Ivoire analyzed. **B)** Frequency of the 59 TE insertions located inside or nearby insecticide-related genes in the Burkina Faso (BF) and Cote d'Ivoire (CI) populations and in our genomes. In red frequencies below 10%. **C)** Frequency of the 438 TE insertions located inside or nearby immune-related genes in the BF and CI populations and in our genomes. In red frequencies below 10%.

Supplemental Table S12. TE insertions in immune-related genes. Genes were obtained from ImmunoDB (Waterhouse *et al.*, 2007). "Promoter" can be YES or NO if found or not in the TE, "dorsal" and "STAT" are the number of binding sites, "# of genomes" is the number of genomes where the TE is found and columns with the genome names can be YES or NO if the insertion is present or not in that particular genome. "Plasmodium responsive" indicates if the gene has been identified as a plasmodium responsive gene according to Ruiz *et al.*, (2019).

Supplemental Table S13. Comparison of the TE copies age distribution (TE landscape) obtained with dnaPipeTE and estimating the identity of the insertions identified by REPET using BLAST. The distributions were compared using a Kolmogorov-Smirnov test and the p-values are reported. Significant differences between superfamilies in each strain are highlighted in yellow.

Supplemental Files

Supplemental File 1. GFF file with the TE annotation for the *AcolN1* genome.

Supplemental File 2. GFF file with the TE annotation for the *DLA112* genome.

Supplemental File 3. GFF file with the TE annotation for the *DLA146* genome.

Supplemental File 4. GFF file with the TE annotation for the *DLA155B* genome.

Supplemental File 5. GFF file with the TE annotation for the *LBV11* genome.

Supplemental File 6. GFF file with the TE annotation for the *LBV88* genome.

Supplemental File 7. GFF file with the TE annotation for the *LBV136* genome.

Supplemental File 8. GFF file with the gene annotation for the *AcolN1* genome.

Supplemental File 9. GFF file with the gene annotation for the *DLA112* genome.

Supplemental File 10. GFF file with the gene annotation for the *DLA146* genome.

Supplemental File 11. GFF file with the gene annotation for the *DLA155B* genome.

Supplemental File 12. GFF file with the gene annotation for the *LBV11* genome.

Supplemental File 13. GFF file with the gene annotation for the *LBV88* genome.

Supplemental File 14. GFF file with the gene annotation for the *LBV136* genome.

Supplemental File 15. TE library for *A. coluzzii* comprising all 520 families in fasta format.

Supplemental File 16. TEs transferred across all genomes in txt format. The first three columns include the chromosome, start and end positions for each insertion are reported. The fourth column the TE family name with a unique number appended after an underscore. Columns five to 11 describe if that insertion was found (PRESENT) or not (ABSENT) or if no information was available (NA).

Supplemental Methods

Sample collection and DNA isolation

We sampled *An. coluzzii* larvae in two cities of Central Africa: Libreville, Gabon, in January 2016 and Douala, Cameroon, in April 2018 (Additional File 1: Table S1). A systematic inspection of potential breeding sites was conducted to determine the presence of *Anopheles* larvae. We manually separated the anopheline from the culicine larvae based on morphological recognition and positioning of their bodies on or under the water surface (Robert, 2017). We collected immature 3rd and 4th stage larvae of *Anopheles* from water bodies using the standard dipping method (Service, 1993). Larvae were stored in 1.5 ml of absolute ethanol. After each daily sampling session, the samples were stored at -20 °C. All the samples were PCR tested to differentiate *An. coluzzii* larvae from *An. gambiae* larvae before library preparation, using primers SINE200_F (TCGCCTTAGACCTTGCGTTA) and SINE200_R (CGCTTCAAGAATTCGAGATAC) (Santolamazza *et al.*, 2008). These primers target a single copy *SINE200* transposable element insertion that is fixed in *An. coluzzii* and absent in *An. gambiae*.

For PacBio sequencing, DNA from a single *An. coluzzii* larva from the *LBV11* site was extracted using the MagAttract HMW DNA extraction kit (Qiagen) following manufacturer's instructions. Briefly, the larva was air-dried and lysed in 240 µl of buffer ATL (proteinase K added) shaking overnight at 56 °C. Next, the DNA was isolated using the MagAttract magnetic beads and eluted twice in 50 µl of buffer AE. The DNA concentration was measured using a Qubit fluorometer. For Nanopore sequencing, DNA from six larvae from each of the five breeding sites was extracted either with the QiaAMP UCP DNA kit (Qiagen) or MagAttract HMW DNA extraction kit (Qiagen). We performed individual larvae extractions as our objective was to use the minimum number of larvae possible to avoid the presence of excess polymorphisms that could affect the genome assembly. For the QiaAMP UCP DNA kit, we followed the manufacturer's instructions. Each larva was air-dried and lysed in 200 µl of buffer AUT (proteinase K added) shaking overnight at 56 °C, then DNA was isolated using a QIAamp UCP MinElute column and eluted twice in 25 µl of buffer AUE. For the MagAttract HMW DNA extraction kit, we followed manufacturer's instructions but using lower buffer amounts to increase DNA concentration. Briefly, each larva was lysed in 120 µl of buffer ATL (proteinase K added) shaking overnight at 56 °C, then DNA was isolated using the MagAttract magnetic beads and eluted twice in 25 µl of

buffer AE. The DNA concentration was measured using a Qubit fluorometer. Both elutions of the same sample were mixed before library preparation. For Illumina sequencing, DNA from one larva from each of the six different breeding sites was extracted following the same extraction protocol as for Nanopore sequencing.

Genome Assemblies

The PacBio sequenced genome was assembled using Canu version 1.8 (Koren *et al.*, 2017) with an estimated genome size of 250Mb and parameters: ‘*stopOnLowCoverage=5, corMinCoverage=0, correctedErrorRate=0.105, CorMhapFilterThreshold=0.0000000002, corMhapOptions="--threshold 0.80 --num-hashes 512 --num-min-matches 3 --ordered-sketch-size 1000 --ordered-kmer-size 14 --min-olap-length 2000 --repeat-idf-scale 50" mhapMemory=60g, mhapBlockSize=500, ovlMerDistinct=0.975*’. The parameter *stopOnLowCoverage* was set to 5 to prevent fragmentation given that some of our samples had medium coverage. *corMinCoverage* was set to 0 to conserve the full length of the reads during the correction stage. *correctedErrorRate* was set to 0.105 following the recommendations in Canu's manual for low coverage genome assemblies. All of the remaining parameters were set to reduce disk space and run time following the recommendations for repetitive genomes. Next, we identified and removed allelic variants using Purge Haplotigs version 1.0.4 (Roach, Schmidt and Borneman, 2018) with the “*-l 15 -m 100 -h 195*” parameters.

The Nanopore genomes were assembled using Canu version 1.8 using the same parameters as previously described, except for *correctedErrorRate* which was set to 0.16, followed by a round of polishing using *racon* version 1.3.3 (Vaser *et al.*, 2017), followed by *nanopolish* version 0.11.1 (Loman, Quick and Simpson, 2015) and *Pilon* version 1.23-0 (Walker *et al.*, 2014) with the *fix* parameter set on ‘*bases*’. *Pilon* requires high coverage short-read data to perform the polishing and these data came from the aforementioned single larvae sequenced from each of the sites. Polishing was not performed for the PacBio sequenced genome given that its assembly had been performed using a single individual and we did not want to incorporate additional variants from the Illumina sequenced individual. In the case of the Nanopore genomes, and although we cannot discard that using Illumina data from an additional individual could introduce novel variants, we did perform the polishing given that those genomes had been obtained using several individuals and polishing increased the BUSCO values. Next, we identified and removed allelic

variants using Purge Haplotigs version 1.0.4 (Roach, Schmidt and Borneman, 2018) with the “*-l 15 -m 100 -h 195*” parameters. Finally, BlobTools version 1.1.1 (Laetsch and Blaxter, 2017) was used to remove contamination from all six genome assemblies taking into consideration fragment sizes, their taxonomic assignment and the coverage using the Illumina reads.

Construction of the curated TE library and *de novo* TE annotation

We ran the TEdenovo pipeline (Flutre *et al.*, 2011) independently on each of the seven genomes with default parameters. The obtained consensus in each genome were further filtered by discarding those generated with (i) only one sequence; (ii) with less than one full-length fragment mapping to the genome; (iii) with less than three full-length copies; and (iv) shorter than 100bp (Additional file 1: Table S2). The remaining consensus were manually curated to remove redundant sequences and artifacts by manual inspection of coverage plots generated using the plotCoverage tool from REPET and visualization of the structural features on the genome browser IGV version 2.4.19 (Robinson *et al.*, 2011).

To ensure that we identified as much of the TE diversity as possible, the TEfam (tefam.biochem.vt.edu) database, which contains the TE libraries for several species of mosquitoes, was used to annotate the seven genomes using RepeatMasker version open-4.0.9 (Smit *et al.* 2015). Because TEfam contains sequences from different species, we used a conservative approach to minimize false positives. We required candidate families to have more than three matches longer than 90% in any genome. Then, the hit with the highest identity from each genome was extracted. These sequences were added to the REPET library and all the consensus were clustered using CD-HIT version 4.8.1 (Fu *et al.*, 2012) with the *-c* and *-s* parameters set to 0.8. These filters ensured that all TEs with an identity greater than 80% throughout more than 80% of their sequence were grouped in the same family. 85 clusters contained sequences only identified by TEfam. The sequences belonging to the same cluster were used to perform a multiple sequence alignment and the consensus were obtained.

The consensus were classified using PASTEC (Hoede *et al.*, 2014) with default parameters. Next their bidirectional best-hits were calculated using BLAST (Camacho *et al.*, 2009) against the TEfam (tefam.biochem.vt.edu), AnoteExcel (Fernández-Medina, Struchiner and Ribeiro,

2011) and Repbase (Bao, Kojima and Kohany, 2015) databases. When more than 80% of a consensus matched to a feature from the databases with an identity higher than 80%, the classification was transferred to the consensus. While not an order *per se*, MITEs were grouped together for subsequent analysis. Additionally, we classified the families based on the conservation of features characteristic of their orders into putative autonomous, putative autonomous lacking terminal inverted repeats (TIRs) or long terminal repeats (LTRs), putative non-autonomous, such as MITEs and TRIMs, and degenerated (Additional file 1: Table S4) (Fonseca *et al* 2019). These classified consensus were used to re-annotate the assembled genomes with the TEannot pipeline using default parameters and we discarded copies whose length overlapped >80% with satellite annotations (Quesneville *et al.*, 2005).

Transfer of TE annotations to the *AcolNI* reference genome

We transferred the euchromatic TE annotations from the six genomes we sequenced to the *AcolNI* genome. First, we built a GFF file composed by the coordinates of two 500 bp long “anchors” adjacent to each TE. We transferred these features considering each pair of anchors as exons from a single gene using the Liftoff tool with the `-exclude_partial -overlap 1 -s 0.8` parameters (Shumate and Salzberg, 2021). We conserved only transfers where both anchors were transferred to the *AcolNI* genome. When both anchors were separated by less than 10 bp we considered the TE to be absent. This allowed us to identify most TE insertions, both present and absent in the *AcolNI* genome. Next, following the same strategy we transferred these regions from the *AcolNI* genome to the other six genomes. This was performed to allow the identification of TEs that were indeed present in these genomes but that had not been initially annotated by REPET. To summarize this information, we built a matrix containing the status for all of the TE insertions transferred to the *AcolNI* genome in every genome. When the anchors were found more than 10 bp away the TE was considered to be present and this was represented with a 1, when the distance was less than 10 bp the TE was considered to be absent, and it was represented with a -1. When any of the anchors was not transferred the TE was considered to be not transferred and represented with a 0. For a graphical description see Figure 1.

Overall, we were able to transfer to the *AcolNI* genome 71.93% to 75.29% of the TEs present in each of the six genomes leading to a total of 67,548 TEs transferred. We then attempt to transferred back these 67,548 from the *AcolNI* reference to the remaining six genomes and we

were able to confidently transfer 53,893 (79.78%) of these regions to at least three of the genomes, with 32,185 (47.65%) transferred to all six genomes. We checked whether the TEs that failed to be transferred were enriched for nested TEs and we found that this was the case: while 68.58% and 51.18% of the TEs located more and less than 500 bp away from other TEs, respectively, were transferred, only 38.26% of the TEs overlapping other TEs were transferred to the six genomes.

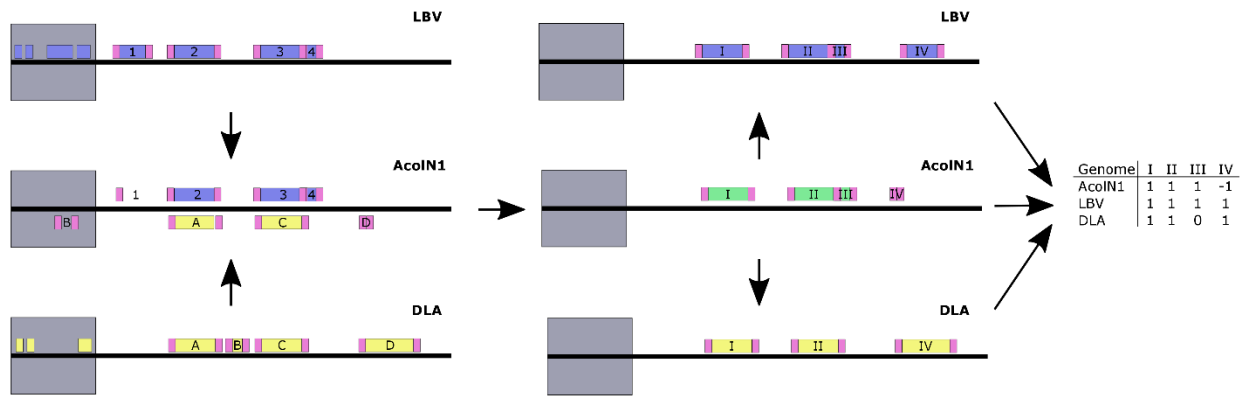


Figure 1. Transfer of TE annotations. Three genomes are shown (LBV, DLA and the reference *AcolN1*) as black lines. The heterochromatin in each genome is displayed using a gray box. The euchromatic TEs from genomes LBV (shown as blue boxes) and DLA (shown as yellow boxes) are transferred to the *AcolN1* genome. This is performed by obtaining 500 bp long anchors from each side of the TE (shown as pink boxes). Then using Liftoff these anchors are transferred to the *AcolN1* genome.

Note that the anchors of TE B were transferred to the heterochromatin and thus this TE is discarded. Also, only one anchor of TE 1 was transferred, therefore it was also discarded. If the distance between both anchors was smaller than 10 (TE D), the TE was considered to be absent in the *AcolN1* genome, however these regions were conserved to assess whether these TEs were present in other genomes. The threshold of 10, instead of a fraction of the size of the TE, was selected to detect solo-LTR insertions given that LTRs represent a small fraction of the total size of the full element. In the case of regions where TEs from two genomes were transferred (for example TE 2 and TE A) if they belonged to the same family the coordinates of the TE with a size closer to the consensus of that family were conserved.

Next the anchors for all TEs that had been transferred to the *AcolNI* genome were transferred to every other genome. Notice that TE IV had not been discovered during the LBV annotation, but it was identified during the transfer. Also, the anchors for TE III were not transferred to the DLA genome. Finally, all this information was summarized in a matrix where 1 represents that the TE was identified, -1 that the TE was absent and 0 when the anchors were not transferred and thus no information of the presence of the TE is available.

Identification of newly described families in other species

We analyzed all 10 available fully sequenced species from the Pyretophorus series, which belongs to the *Cellia* subgenus. We also included an additional five *Anopheles* species, three from each of the other series from the *Cellia* subgenus and two from the other subgenera with available fully sequenced species. As outgroups we included the genomes of *Cx.*

quinquefasciatus, *Ae. aegypti* and *D. melanogaster*. RepeatMasker version open-4.0.9 (Smit *et al.* 2015) was run with default parameters using the 64 newly described families as the library on the following genomes: *An. albimanus* (AalbS2), *An. atroparvus* (AatrE3), *An. farauti* (AfarF2), *An. funestus* (AfunF3), *An. stephensi* (AsteS1), *An. epiroticus* (AepiE1), *An. christyi* (AchrA1), *An. merus* (AmerM2), *An. gambiae* (AgamP4), *An. coluzzii* (*AcolNI*), *An. melas* (AmelC2), *An. arabiensis* (AaraD1), *An. quadriannulatus* (AquaS1), *An. bwambae* (Abwa2) and *An. fontenillei* (ASM881789v1), *Cx. quinquefasciatus* (CulPip1.0), *Ae. aegypti* (AaegL5.0) and *D. melanogaster* (ISO1 release 6).

Insertion frequency estimation in rural populations

We used the PoPoolationTE2 v-1.10.03 pipeline (Kofler, Gómez-Sánchez and Schlötterer, 2016) to compute the TE insertion frequencies in the Ivory Coast and Burkina Faso populations from the Ag1000G project (*Anopheles gambiae* 1000 Genomes Consortium, 2020).

The TE hierarchy and the TE-merged reference files were created using the TE library previously created and the *AcolNI* reference genome. We used RepeatMasker version open-4.0.9 (Smit, Hubley and Green, 2013-2015) to mask the TE sequences in the reference and concatenated the masked reference file with the full TE library fasta. Next, we mapped the sequencing data from each sample in fastq format against the concatenated reference using bwa version 0.7.12-r1039 (Li and Durbin, 2009) with default parameters. The obtained files were

converted to bam format, sorted and indexed using samtools version 1.9-55-g8a4d196 (Li *et al.*, 2009).

Using the ppileup tool from PopoolationTE2 we generated a pileup from all the bam files. Next, the signatures were identified with the identifySignatures tool using the separate mode and requiring a minimum count of 3. The output file was processed with the frequency tool and the pairupSignatures tool using a minimum distance of 200 and a maximum distance of 300 as recommended in the manual. Once we obtained the insertion locations we converted them to GenomicRanges objects ignoring strand information and used the findOverlaps function from the GenomicAlignments R package (Lawrence *et al.*, 2013) to identify overlaps with the previously identified TEs of interest. We only considered matches where the overlapping TEs belonged to the same superfamily.

Supplemental References

- Bao, W., Kojima, K. K. and Kohany, O. (2015) 'Repbase Update, a database of repetitive elements in eukaryotic genomes', *Mobile DNA*, 6(1), pp. 4-9.
- Camacho, C., Coulouris, G., Avagyan, V., Ma, N., Papadopoulos, J., Bealer, K. and Madden, T. L. (2009) 'BLAST+: architecture and applications', *BMC bioinformatics*, 10, pp. 421-421.
- Anopheles gambiae 1000 Genomes Consortium (2020) 'Genome variation and population structure among 1,142 mosquitoes of the African malaria vector species *Anopheles gambiae* and *Anopheles coluzzii*', *Genome Research*, 30(10):1533-1546.
- Fernández-Medina, R. D., Struchiner, C. J. and Ribeiro, J. M. C. (2011) 'Novel transposable elements from *Anopheles gambiae*', *BMC Genomics*, 12(1), pp. 260-260.
- Flutre, T., Duprat, E., Feuillet, C. and Quesneville, H. (2011) 'Considering Transposable Element Diversification in De Novo Annotation Approaches', *PLOS ONE*, 6(1), pp. e16526.
- Fu, L., Niu, B., Zhu, Z., Wu, S. and Li, W. (2012) 'CD-HIT: accelerated for clustering the next-generation sequencing data', *Bioinformatics (Oxford, England)*, 28(23), pp. 3150-3152.
- Hoede, C., Arnoux, S., Moisset, M., Chaumier, T., Inizan, O., Jamilloux, V. and Quesneville, H. (2014) 'PASTEC: An Automatic Transposable Element Classification Tool', *PLOS ONE*, 9(5), pp. e91929.
- Kofler, R., Gómez-Sánchez, D. and Schlötterer, C. (2016) 'PoPoolationTE2: Comparative Population Genomics of Transposable Elements Using Pool-Seq', *Molecular biology and evolution*, 33(10), pp. 2759-2764.
- Koren, S., Walenz, B. P., Berlin, K., Miller, J. R., Bergman, N. H. and Phillippy, A. M. (2017) 'Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation', *Genome research*, 27(5), pp. 722-736.
- Laetsch, D. R. and Blaxter, M. L. (2017) 'BlobTools: Interrogation of genome assemblies', *F1000Research*, 6(1287).
- Lawrence, M., Huber, W., Pagès, H., Aboyoun, P., Carlson, M., Gentleman, R., Morgan, M. T. and Carey, V. J. (2013) 'Software for Computing and Annotating Genomic Ranges', *PLOS Computational Biology*, 9(8), pp. e1003118.

- Li, H. and Durbin, R. (2009) 'Fast and accurate short read alignment with Burrows–Wheeler transform', *Bioinformatics*, 25(14), pp. 1754-1760.
- Li, H., Handsaker, B., Wysoker, A., Fennell, T., Ruan, J., Homer, N., Marth, G., Abecasis, G., Durbin, R. and Genome Project Data Processing, S. (2009) 'The Sequence Alignment/Map format and SAMtools', *Bioinformatics*, 25(16), pp. 2078-2079.
- Loman, N. J., Quick, J. and Simpson, J. T. (2015) 'A complete bacterial genome assembled de novo using only nanopore sequencing data', *Nature Methods*, 12(8), pp. 733-735.
- Quesneville, H., Bergman, C. M., Andrieu, O., Autard, D., Nouaud, D., Ashburner, M. and Anxolabehere, D. (2005) 'Combined evidence annotation of transposable elements in genome sequences', *PLoS computational biology*, 1(2), pp. 166-175.
- Roach, M. J., Schmidt, S. A. and Borneman, A. R. (2018) 'Purge Haplotigs: allelic contig reassignment for third-gen diploid genome assemblies', *BMC Bioinformatics*, 19(1), pp. 460.
- Robinson, J. T., Thorvaldsdóttir, H., Winckler, W., Guttman, M., Lander, E. S., Getz, G. and Mesirov, J. P. (2011) 'Integrative genomics viewer', *Nature biotechnology*, 29(1), pp. 24-26.
- Santolamazza, F., Mancini, E., Simard, F., Qi, Y., Tu, Z. and Della Torre, A. (2008) 'Insertion polymorphisms of SINE200 retrotransposons within speciation islands of *Anopheles gambiae* molecular forms', *Malaria Journal*, 7, pp. 1-10.
- Shumate, A. and Salzberg, S. L. (2021) 'Liftoff: accurate mapping of gene annotations', *Bioinformatics*, 37(12), pp. 1639–1643.
- Smit, A., Hubley, R. and Green, P. (2013-2015) *RepeatMasker Open-4.0*. Available at: <http://www.repeatmasker.org>.
- Vaser, R., Sović, I., Nagarajan, N. and Šikić, M. (2017) 'Fast and accurate de novo genome assembly from long uncorrected reads', *Genome research*, 27(5), pp. 737-746.
- Walker, B. J., Abeel, T., Shea, T., Priest, M., Abouelliel, A., Sakthikumar, S., Cuomo, C. A., Zeng, Q., Wortman, J., Young, S. K. and Earl, A. M. (2014) 'Pilon: An Integrated Tool for Comprehensive Microbial Variant Detection and Genome Assembly Improvement', *PLOS ONE*, 9(11), pp. e112963.